
Entropic Gromov-Wasserstein Distances: Stability and Algorithms

Gabriel Rioux

Center for Applied Mathematics
Cornell University
Ithaca, NY 14853.
ger84@cornell.edu

Ziv Goldfeld

Department of Electrical and Computer Engineering
Cornell University
Ithaca, NY 14853.
goldfeld@cornell.edu

Kengo Kato

Department of Statistics and Data Science
Cornell University
Ithaca, NY 14853.
kk976@cornell.edu

Abstract

The Gromov-Wasserstein (GW) distance quantifies discrepancy between metric measure spaces, but suffers from computational hardness. The entropic Gromov-Wasserstein (EGW) distance serves as a computationally efficient proxy for the GW distance. Recently, it was shown that the quadratic GW and EGW distances admit variational forms that tie them to the well-understood optimal transport (OT) and entropic OT (EOT) problems. By leveraging this connection, we establish convexity and smoothness properties of the objective in this variational problem. This results in the first efficient algorithms for solving the EGW problem that are subject to formal guarantees in both the convex and non-convex regimes.

1 Introduction

The Gromov-Wasserstein (GW) distance compares probability distributions that are supported on possibly distinct metric spaces by aligning them with one another. Given two metric measure (mm) spaces $(\mathcal{X}_0, d_0, \mu_0)$ and $(\mathcal{X}_1, d_1, \mu_1)$, the (p, q) -GW distance between them is

$$D_{p,q}(\mu_0, \mu_1) := \inf_{\pi \in \Pi(\mu_0, \mu_1)} \left(\int |d_0^q(x, x') - d_1^q(y, y')|^p d\pi \otimes \pi(x, y, x', y') \right)^{\frac{1}{p}}, \quad (1)$$

where $\Pi(\mu_0, \mu_1)$ is the set of couplings between μ_0 and μ_1 . This approach, proposed in [25], is an optimal transport (OT) based L^p relaxation of the classical Gromov-Hausdorff distance between metric spaces. The GW distance defines a metric on the quotient space of all mm spaces modulo obtained by identifying isomorphic mm spaces (i.e. the underlying measures μ_0, μ_1) are such that $\mu_0 \circ T^{-1} = \mu_1$ for some isometry $T : \mathcal{X}_0 \rightarrow \mathcal{X}_1$). From an applied standpoint, a solution to the GW problem between two heterogeneous datasets yields both a quantification of discrepancy, and an optimal matching π^* between them. As such, the GW distance has seen many applications, encompassing single-cell genomics [5, 15], alignment of language models [1], shape matching [23, 24], graph matching [39, 40], heterogeneous domain adaptation [41], and generative modeling [8].

Exact computation of the GW distance is a quadratic assignment problem, which is known to be NP-complete [11]. The computational intractability of the GW problem in (1) has inspired several reformulations that aim to alleviate this issue. Recent approaches include slicing [38], relaxing

the strict marginal constraints using f -divergence penalties [33], and optimizing over bi-directional maps [44]. While these methods offer certain advantages, it is the approach based on entropic regularization [29, 36] that is most frequently used in application. In [29], it is proposed to solve the entropic Gromov-Wasserstein problem (EGW) via a mirror descent algorithm with a complexity of $O(N^3)$ for marginals supported on N distinct points (see, e.g., Remark 1 in [29]). The follow-up work [32] proposes a low-rank variant of the EGW problem which can be solved in linear time, wherein only couplings admitting a certain low-rank structure are considered. As an intermediate step of their analysis, they show that the complexity of mirror descent can be reduced to $O(N^2)$ by assuming that the matrices of pairwise costs admit a low-rank decomposition (without imposing any structure on the couplings). This decomposition holds, for instance, when the cost is quadratic and N dominates the ambient dimensions. Although mirror descent seems to solve the EGW problem well in practice, formal guarantees concerning convergence rates or local optimality are lacking.

The goal of this work is to address the computational gap described above, targeting algorithms with non-asymptotic guarantees and establishing convexity regimes of the EGW problem—all of which are consequences of a new stability analysis of the EGW variational representation from [43].

2 Notation and preliminaries

For a topological space S , we let $\mathcal{P}(S)$ be the collection of all Borel probability distributions on S . The Frobenius inner product on $\mathbb{R}^{d_0 \times d_1}$ is defined by $\langle \mathbf{A}, \mathbf{B} \rangle_F = \text{tr}(\mathbf{A}^\top \mathbf{B})$; the associated norm is denoted by $\|\cdot\|_F$. A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is ρ -weakly convex if $f + \frac{\rho}{2}\|\cdot\|^2$ is convex, it is L -smooth if its gradient is L -Lipschitz. For a Fréchet differentiable map $F: \bar{U} \rightarrow V$ between normed vector spaces U and V , we denote the derivative of F at the point $u \in U$ evaluated at $v \in V$ by $DF_{[u]}(v)$. We adopt the shorthands $a \wedge b = \min\{a, b\}$ and $a \vee b = \max\{a, b\}$.

2.1 Entropic optimal transport

Entropic regularization transforms the linear OT problem into a strongly convex one. Given distributions $\mu_i \in \mathcal{P}(\mathbb{R}^{d_i})$, $i = 0, 1$, and a cost function $c: \mathbb{R}^{d_0} \times \mathbb{R}^{d_1} \rightarrow \mathbb{R}$, the primal EOT problem is obtained by regularizing the standard OT problem via the Kullback-Leibler (KL) divergence, $\text{OT}_\varepsilon(\mu_0, \mu_1) = \inf_{\pi \in \Pi(\mu_0, \mu_1)} \int c d\pi + \varepsilon \text{D}_{\text{KL}}(\pi \| \mu_0 \otimes \mu_1)$, where $\varepsilon > 0$ is a regularization parameter and $\text{D}_{\text{KL}}(\mu \| \nu) = \int \log\left(\frac{d\mu}{d\nu}\right) d\mu$, if $\mu \ll \nu$, and ∞ , otherwise. Classical OT is obtained from the above by setting $\varepsilon = 0$. When $c \in L^1(\mu_0 \otimes \mu_1)$, EOT admits the following dual formulation,

$$\text{OT}_\varepsilon(\mu_0, \mu_1) = \sup_{(\varphi_0, \varphi_1) \in L^1(\mu_0) \times L^1(\mu_1)} \int \varphi_0 d\mu_0 + \int \varphi_1 d\mu_1 - \varepsilon \int e^{\frac{\varphi_0 \oplus \varphi_1 - c}{\varepsilon}} d\mu_0 \otimes \mu_1 + \varepsilon,$$

where $\varphi_0 \oplus \varphi_1(x, y) = \varphi_0(x) + \varphi_1(y)$. The set of solutions to the dual problem coincides with the set of solutions to the so-called Schrödinger system,

$$\int e^{\frac{\varphi_0(x) + \varphi_1(y) - c(x, y)}{\varepsilon}} d\mu_1 = 1, \quad \mu_0\text{-a.e. } x \in \mathbb{R}^{d_0}, \quad \int e^{\frac{\varphi_0(x) + \varphi_1(y) - c(x, y)}{\varepsilon}} d\mu_0 = 1, \quad \mu_1\text{-a.e. } y \in \mathbb{R}^{d_1}, \quad (2)$$

for $(\varphi_0, \varphi_1) \in L^1(\mu_0) \times L^1(\mu_1)$. A pair $(\varphi_0, \varphi_1) \in L^1(\mu_0) \times L^1(\mu_1)$ solving (2) is known to be a.s. unique up to additive constants in the sense that any other solution $(\bar{\varphi}_0, \bar{\varphi}_1)$ satisfies $\bar{\varphi}_0 = \varphi_0 + a$ μ_0 -a.s. and $\bar{\varphi}_1 = \varphi_1 - a$ μ_1 -a.s. for some $a \in \mathbb{R}$. The unique EOT coupling π_ε is characterized by $\frac{d\pi_\varepsilon}{d\mu_0 \otimes \mu_1}(x, y) = e^{\frac{\varphi_0(x) + \varphi_1(y) - c(x, y)}{\varepsilon}}$, and, under some additional conditions on the cost and marginals, (2) admits a pair of continuous solutions which is unique up to additive constants and satisfies the system everywhere, i.e., at all points $(x, y) \in \mathbb{R}^{d_0} \times \mathbb{R}^{d_1}$. We call such continuous solutions EOT potentials. The reader is referred to [28] for a comprehensive overview of EOT.

2.2 Entropic Gromov-Wasserstein distance

This work studies stability and computational aspects of the entropically regularized GW distance under the quadratic and the inner product cost. By analogy to OT, EGW serves as a proxy of the standard (p, q) -GW distance. From here on out we instantiate the mm spaces as the Euclidean spaces $(\mathbb{R}^{d_i}, \|\cdot\|, \mu_i)$, for $i = 0, 1$, and proceed to define the EGW distance for the quadratic cost.

The quadratic EGW distance, which corresponds to the $p = q = 2$ case, is defined as

$$S_\varepsilon(\mu_0, \mu_1) = \inf_{\pi \in \Pi(\mu_0, \mu_1)} \int \left(\|x - x'\|^2 - \|y - y'\|^2 \right)^2 d\pi \otimes \pi(x, y, x', y') + \varepsilon D_{\text{KL}}(\pi \| \mu_0 \otimes \mu_1). \quad (3)$$

One readily verifies that, like the standard GW distance, EGW is invariant to isometric actions on the marginal spaces such as orthogonal rotations and translations. When μ_0, μ_1 are centered, which we may assume without loss of generality, the EGW distance decomposes as

$$\begin{aligned} S_\varepsilon(\mu_0, \mu_1) &= S_1(\mu_0, \mu_1) + S_{2,\varepsilon}(\mu_0, \mu_1), \\ S_1(\mu_0, \mu_1) &= \int \|x - x'\|^4 d\mu_0 \otimes \mu_0(x, x') + \int \|y - y'\|^4 d\mu_1 \otimes \mu_1(y, y') - 4M_2(\mu_0)M_2(\mu_1), \quad (4) \\ S_{2,\varepsilon}(\mu_0, \mu_1) &= \inf_{\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}} 32\|\mathbf{A}\|_F^2 + \text{OT}_{\mathbf{A},\varepsilon}(\mu_0, \mu_1), \end{aligned}$$

where $\text{OT}_{\mathbf{A},\varepsilon}(\mu_0, \mu_1)$ is the EOT problem with the cost function $c_{\mathbf{A}} : (x, y) \in \mathbb{R}^{d_0} \times \mathbb{R}^{d_1} \mapsto -4\|x\|^2\|y\|^2 - 32x^\top \mathbf{A}y$ and regularization parameter ε . Moreover, the infimum is achieved at some $\mathbf{A}^* \in \mathcal{D}_M := [-M/2, M/2]^{d_0 \times d_1}$ for any $M \geq \sqrt{M_2(\mu_0)M_2(\mu_1)} =: M_{\mu_0, \mu_1}$. The proof of Theorem 1 in [43] demonstrates that if μ_0 and μ_1 are centered and π_* is optimal for the original EGW formulation, then $\mathbf{A}^* = \frac{1}{2} \int xy^\top d\pi_*(x, y)$ is optimal for $S_{2,\varepsilon}$ and $\pi_* = \pi_{\mathbf{A}^*}$, where $\pi_{\mathbf{A}^*}$ is the unique EOT coupling for $\text{OT}_{\mathbf{A}^*,\varepsilon}(\mu_0, \mu_1)$. Corollary 1 ahead expands on this connection by establishing a one-to-one correspondence between solutions of S_ε and $S_{2,\varepsilon}$.

Although (4) illustrates a connection between the EGW and EOT problems, the outer minimization over \mathcal{D}_M necessitates studying EOT with an *a priori* unknown cost function $c_{\mathbf{A}}$.

A similar decomposition holds for the inner product GW problem, where the difference of squared Euclidean norms is replaced by a difference of inner products. In that case, $F_\varepsilon(\mu_0, \mu_1) = F_1(\mu_0, \mu_1) + F_{2,\varepsilon}(\mu_0, \mu_1)$, for $F_1(\mu_0, \mu_1) = \int |\langle x, x' \rangle|^2 d\mu_0 \otimes \mu_0(x, x') + \int |\langle y, y' \rangle|^2 d\mu_0 \otimes \mu_0(y, y')$, and $F_{2,\varepsilon}(\mu_0, \mu_1) = \inf_{\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}} 8\|\mathbf{A}\|_F^2 + \text{IOT}_{\mathbf{A},\varepsilon}(\mu_0, \mu_1)$, with the distinction that no centering is needed and $\text{IOT}_{\mathbf{A},\varepsilon}(\mu_0, \mu_1)$ is the EOT problem with the cost function $c_{\mathbf{A}}(x, y) = -8x^\top \mathbf{A}y$. We restrict our attention to the quadratic EGW problem, similar results hold in the inner product case.

3 Stability of entropic Gromov-Wasserstein distances

We now analyze the stability of the EGW problem with respect to the matrix \mathbf{A} appearing in its variational form (4). Specifically, we characterize the first and second derivatives of the objective function whose optimization defines $S_{2,\varepsilon}$ which elucidates its convexity properties and enables us to devise novel approaches for computing the EGW distance with formal convergence guarantees. Throughout this section, we restrict attention to compactly supported distributions, as some of the technical details do not directly translate to the unbounded setting (e.g., the proof of Lemma 2).

Fix compactly supported distributions $(\mu_0, \mu_1) \in \mathcal{P}(\mathbb{R}^{d_0}) \times \mathcal{P}(\mathbb{R}^{d_1})$ and some $\varepsilon > 0$. Let

$$\Phi : \mathbf{A} \in \mathbb{R}^{d_0 \times d_1} \mapsto 32\|\mathbf{A}\|_F^2 + \text{OT}_{\mathbf{A},\varepsilon}(\mu_0, \mu_1)$$

denote the objective in $S_{2,\varepsilon}(\mu_0, \mu_1)$. We first characterize the derivatives of Φ and then prove that this map is weakly convex and L -smooth.

Proposition 1 (First and second derivatives). $\Phi : \mathbf{A} \in \mathbb{R}^{d_0 \times d_1} \mapsto 32\|\mathbf{A}\|_F^2 + \text{OT}_{\mathbf{A},\varepsilon}(\mu_0, \mu_1)$ is smooth, coercive, and has first and second-order Fréchet derivatives at $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$ given by

$$D\Phi_{[\mathbf{A}]}(\mathbf{B}) = 64 \text{tr}(\mathbf{A}^\top \mathbf{B}) - 32 \int x^\top \mathbf{B}y d\pi_{\mathbf{A}}(x, y),$$

$$D^2\Phi_{[\mathbf{A}]}(\mathbf{B}, \mathbf{C}) = 64 \text{tr}(\mathbf{B}^\top \mathbf{C}) + 32\varepsilon^{-1} \int x^\top \mathbf{B}y \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) - 32x^\top \mathbf{C}y \right) d\pi_{\mathbf{A}}(x, y),$$

where $\mathbf{B}, \mathbf{C} \in \mathbb{R}^{d_0 \times d_1}$, $\pi_{\mathbf{A}}$ is the unique EOT coupling for $\text{OT}_{\mathbf{A},\varepsilon}(\mu_0, \mu_1)$, and $(h_0^{\mathbf{A},\mathbf{C}}, h_1^{\mathbf{A},\mathbf{C}})$ is the unique (up to additive constants) pair of functions in $\mathcal{C}(\text{spt}(\mu_0)) \times \mathcal{C}(\text{spt}(\mu_1))$ satisfying

$$\begin{aligned} \int \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) - 32x^\top \mathbf{C}y \right) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x,y)}{\varepsilon}} d\mu_1(y) &= 0, \quad \forall x \in \text{spt}(\mu_0), \\ \int \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) - 32x^\top \mathbf{C}y \right) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x,y)}{\varepsilon}} d\mu_0(x) &= 0, \quad \forall y \in \text{spt}(\mu_1). \end{aligned} \quad (5)$$

Here, $(\varphi_0^A, \varphi_1^A)$ is any pair of EOT potentials for $\text{OT}_{A,\varepsilon}(\mu_0, \mu_1)$.

Proposition 1 essentially follows from the implicit mapping theorem, which enables us to compute the Fréchet derivative of the EOT potentials for $\text{OT}_{(\cdot),\varepsilon}(\mu_0, \mu_1)$ using the Schrödinger system (2). The derivative of $\text{OT}_{(\cdot),\varepsilon}(\mu_0, \mu_1)$, which is a simple function of the EOT potentials, is then readily obtained. By differentiating the Frobenius norm, this yields the derivative of Φ . See Appendix A.1.

As $D\Phi_{[A]}(\mathbf{B}) = \langle 64\mathbf{A} - 32 \int xy^\top d\pi_{\mathbf{A}}(x, y), \mathbf{B} \rangle_F$, we can interpret $64\mathbf{A} - 32 \int xy^\top d\pi_{\mathbf{A}}(x, y)$ as the gradient of Φ which we denote $D\Phi_{[A]}$. This perspective is utilized in Section 4 when studying computational guarantees for the EGW distance, as it is simpler to view iterates as matrices.

As a direct corollary to Proposition 1, we provide an (implicit) characterization of the stationary points of Φ and connect its minimizers to solutions of S_ε . Details are provided in Appendix A.2.

Corollary 1 (Stationary points and correspondence between S_ε and $S_{2,\varepsilon}$).

- (i) A matrix $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$ is a stationary point of Φ if and only if $\mathbf{A} = \frac{1}{2} \int xy^\top d\pi_{\mathbf{A}}(x, y)$. As Φ is coercive, all minimizers of Φ are stationary points and hence contained in $\mathcal{D}_{M_{\mu_0, \mu_1}}$.
- (ii) If μ_0 and μ_1 are centered, then a given matrix \mathbf{A} minimizes Φ if and only if $\pi_{\mathbf{A}}$ is optimal for S_ε and satisfies $\frac{1}{2} \int xy^\top d\pi_{\mathbf{A}}(x, y) = \mathbf{A}$.
- (iii) Suppose μ_0 and μ_1 are centered. If S_ε admits a unique optimal coupling π_* , then Φ admits a unique minimizer \mathbf{A}^* and $\pi_* = \pi_{\mathbf{A}^*}$. Conversely, if Φ admits a unique minimizer \mathbf{A}^* , then $\pi_{\mathbf{A}^*}$ is a unique optimal coupling for S_ε .

Although the second derivative of Φ involves the implicitly defined functions $(h_0^{A,C}, h_1^{A,C})$, its maximal and minimal eigenvalues, $\lambda_{\max}(D^2\Phi_{[A]})$ and $\lambda_{\min}(D^2\Phi_{[A]})$, can be controlled which enables us to characterize convexity and smoothness of Φ .

Theorem 1 (Convexity and L -smoothness). *The map Φ is weakly convex with parameter at most $32^2\varepsilon^{-1}\sqrt{M_4(\mu_0)M_4(\mu_1)} - 64$ and, if $\sqrt{M_4(\mu_0)M_4(\mu_1)} < \frac{\varepsilon}{16}$, then it is strictly convex and admits a unique minimizer. Moreover, for any $M > 0$, Φ is L -smooth on \mathcal{D}_M with $L \leq 64 \vee \left(32^2\varepsilon^{-1}\sqrt{M_4(\mu_0)M_4(\mu_1)} - 64\right)$.*

Theorem 1 follows from Proposition 1 by considering the variational form of the maximal and minimal eigenvalues; see Appendix A.3 for details. In general, optimal EGW couplings may not be unique. Theorem 1 provides sufficient conditions for uniqueness of solutions to both $S_{2,\varepsilon}$ and the EGW problem by the connection discussed in Corollary 1 when the marginals are centered.

4 Computational guarantees

Building on this stability theory, we now study computation of the EGW problem. The goal is to compute the distance between two discrete distributions $\mu_0 \in \mathcal{P}(\mathbb{R}^{d_0})$ and $\mu_1 \in \mathcal{P}(\mathbb{R}^{d_1})$ supported on N_0 and N_1 atoms $(x^{(i)})_{i=1}^{N_0}$ and $(y^{(j)})_{j=1}^{N_1}$, respectively. In light of the decomposition (4), we focus on $S_{2,\varepsilon}$, which is given by a smooth optimization problem whose convexity depends on the value of ε . Throughout, we treat $D\Phi_{[A]}$, for $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$, as the matrix $64\mathbf{A} - 32 \int xy^\top d\pi_{\mathbf{A}}(x, y)$.

4.1 Inexact Oracle Methods

As these problems are already d_0d_1 -dimensional and computing the second Fréchet derivative of Φ may be infeasible (in particular, it requires solving Eq. (5)), we focus on first-order methods. Given the regularity of the $S_{2,\varepsilon}$ optimization problem, standard out-of-the-box numerical routines are likely to yield good results in practice. However, to provide meaningful formal guarantees one must account for the fact that evaluation of Φ and its gradient requires computing the corresponding EOT plan, which entails an approximation. We model this under the scope of gradient methods with inexact gradient oracles [13, 16, 17].

For a fixed $\varepsilon > 0$ and μ_0, μ_1 as above, we seek to solve $\min_{\mathbf{A} \in \mathcal{D}_M} 32\|\mathbf{A}\|_F^2 + \text{OT}_{A,\varepsilon}(\mu_0, \mu_1)$, where $M > M_{\mu_0, \mu_1}$, which guarantees that all the optimizers are within the optimization domain (cf. Corollary 1). As we are in the discrete setting, the EOT coupling $\pi^{\mathbf{A}}$ for $\text{OT}_{A,\varepsilon}(\mu_0, \mu_1)$, $\mathbf{A} \in \mathcal{D}_M$,

is represented by $\Pi^{\mathbf{A}} \in \mathbb{R}^{N_0 \times N_1}$, where $\Pi_{ij}^{\mathbf{A}} = \pi^{\mathbf{A}}(x^{(i)}, y^{(j)})$. The inexact oracle paradigm assumes that, for any $\mathbf{A} \in \mathcal{D}_M$, we have access to a δ -oracle $\tilde{\Pi}^{\mathbf{A}}$ for $\Pi^{\mathbf{A}}$ with $\|\tilde{\Pi}^{\mathbf{A}} - \Pi^{\mathbf{A}}\|_{\infty} < \delta$. Such oracles can be obtained, for instance, by Sinkhorn's algorithm [12, 35].

Proposition 2 (Inexact oracle via Sinkhorn iterations). *Fix $\delta > 0$. Then, Sinkhorn's algorithm (Algorithm 3) returns a δ -oracle approximation $\tilde{\Pi}^{\mathbf{A}}$ of $\Pi^{\mathbf{A}}$ in at most \tilde{k} iterations, where \tilde{k} depends only on $\mu_0, \mu_1, \mathbf{A}, \delta$, and ε , and is given explicitly in (17).*

The proof of Proposition 2 follows by combining a number of known results, see Appendix C. With these preparations, we first discuss the case where Φ is known to be convex on \mathcal{D}_M .

4.2 Convex case

Assume that Φ is convex on \mathcal{D}_M , e.g., under the setting of Theorem 1. As convexity implies that the minimal eigenvalue of $D^2\Phi_{[\mathbf{A}]}$ is positive for any $\mathbf{A} \in \mathcal{D}_M$, Theorem 1 further yields that Φ is 64-smooth. With that, we can apply inexact oracle first-order method from [13]. To describe the approach, assume that we are given a δ -oracle $\tilde{\Pi}^{\mathbf{A}}$ for the EOT plan $\Pi^{\mathbf{A}}$ for $\text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1)$, and define the corresponding gradient approximation

$$\tilde{D}\Phi_{[\mathbf{A}]} = 64\mathbf{A} - 32 \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} x^{(i)}(y^{(j)})^{\top} \tilde{\Pi}_{ij}^{\mathbf{A}}. \quad (6)$$

Algorithm 1 Fast gradient method with inexact oracle

- Fix $L = 64$ and let $\alpha_k = \frac{k+1}{2}$, and $\tau_k = \frac{2}{k+3}$
- 1: $k \leftarrow 0, \mathbf{A}_0 \leftarrow \mathbf{0}, \mathbf{G}_0 \leftarrow \tilde{D}\Phi_{[\mathbf{A}_0]}, \mathbf{W}_0 \leftarrow \alpha_0 \mathbf{G}_0$
 - 2: **while** stopping condition is not met **do**
 - 3: $\mathbf{D}_k \leftarrow \mathbf{A}_k - L^{-1} \mathbf{G}_k$
 - 4: $\mathbf{B}_k \leftarrow \frac{M}{2} \text{sign}(\mathbf{D}_k) \min\left(\frac{2}{M} \|\mathbf{D}_k\|, 1\right)$
 - 5: $\mathbf{C}_k \leftarrow \frac{M}{2} \text{sign}\left(-\frac{\mathbf{W}_k}{L}\right) \min\left(\frac{2}{M} \|\frac{\mathbf{W}_k}{L}\|, 1\right)$
 - 6: $\mathbf{A}_{k+1} \leftarrow \tau_k \mathbf{C}_k + (1 - \tau_k) \mathbf{B}_k$
 - 7: $\mathbf{G}_{k+1} \leftarrow \tilde{D}\Phi_{[\mathbf{A}_{k+1}]}$
 - 8: $\mathbf{W}_{k+1} \leftarrow \mathbf{W}_k + \alpha_{k+1} \mathbf{G}_{k+1}$
 - 9: $k \leftarrow k + 1$
 - 10: **return** \mathbf{B}_k
-

We now present the algorithm and follow it with formal convergence guarantees.

The sign, min, and multiplication operations in Algorithm 1 are applied entrywise. Due to inexactness, stopping conditions based on insufficient progress of functions values or setting a threshold on the norm of the gradient require care. A condition based on the number of iterations is discussed in Remark 1.

We now provide formal convergence guarantees for Algorithm 1.

Theorem 2 (Fast convergence rates). *Assume that Φ is convex and L -smooth on \mathcal{D}_M and that $\tilde{\Pi}^{\mathbf{A}}$ is a δ -oracle for $\Pi^{\mathbf{A}}$. Then, the iterates \mathbf{B}_k in Algorithm 1 with $\tilde{D}\Phi_{[\mathbf{A}_k]}$ given by (6) satisfy $\Phi(\mathbf{B}_k) - \Phi(\mathbf{B}^*) \leq \frac{2L\|\mathbf{B}^*\|_F^2}{(k+1)(k+2)} + 3\delta'$, where \mathbf{B}^* is a global minimizer of Φ and $\delta' = 32M\delta \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} \|x^{(i)}(y^{(j)})^{\top}\|_1$ where $\|\cdot\|_1$ denotes the entrywise 1-norm. Moreover, for any $\eta >$*

$3\delta'$, Algorithm 1 requires at most $k = \left\lceil -\frac{3}{2} + \frac{1}{2} \sqrt{1 + \frac{8L\|\mathbf{B}^\|_F^2}{\eta - 3\delta'}} \right\rceil \leq \left\lceil -\frac{3}{2} + \frac{1}{2} \sqrt{1 + \frac{128M^2 d_0^2 d_1^2}{\eta - 3\delta'}} \right\rceil$ iterations to achieve an η -approximate solution.*

The proof of Theorem 2, given in Appendix A.4, follows from Theorem 2.2 in [13] after casting our problem as an instance of their setting. Some implications of Theorem 2 are discussed next.

Remark 1 (Optimal rates and stopping conditions). *First, consider the convergence rate of the function values. The first term on the right-hand side exhibits the optimal complexity bound for smooth constrained optimization of $O(1/k^2)$ (cf., e.g., [27]). The second term accounts for the underlying oracle error. Notably, the progress of the optimization procedure and the oracle error are completely decoupled in this bound.*

Next, observe that all terms involved in the upper bound for the number of iterations are explicit as soon as a desired precision η is chosen since the oracle error δ can be fixed according to Proposition 2. Consequently, it can be used as an explicit stopping condition for Algorithm 1.

4.3 General case

Algorithm 2 Adaptive gradient method with inexact oracle

Given $\mathbf{C}_0 \in \mathcal{D}_M$, fix the sequences $\beta_k = \frac{1}{2L}$, $\gamma_k = \frac{k}{4L}$, and $\tau_k = \frac{2}{k+2}$.

- 1: $k \leftarrow 1$, $\mathbf{A}_1 \leftarrow \mathbf{C}_0$, $\mathbf{G}_1 \leftarrow \tilde{D}\Phi_{[\mathbf{A}_1]}$
- 2: **while** stopping condition is not met **do**
- 3: $\mathbf{D}_k \leftarrow \mathbf{A}_k - \beta_k \mathbf{G}_k$
- 4: $\mathbf{B}_k \leftarrow \frac{M}{2} \text{sign}(\mathbf{D}_k) \min\left(\frac{2}{M} \|\mathbf{D}_k\|, 1\right)$
- 5: $\mathbf{E}_k \leftarrow \mathbf{C}_{k-1} - \gamma_k \mathbf{G}_k$
- 6: $\mathbf{C}_k \leftarrow \frac{M}{2} \text{sign}(\mathbf{E}_k) \min\left(\frac{2}{M} \|\mathbf{E}_k\|, 1\right)$
- 7: $\mathbf{A}_{k+1} \leftarrow \tau_k \mathbf{C}_k + (1 - \tau_k) \mathbf{B}_k$
- 8: $\mathbf{G}_{k+1} \leftarrow \tilde{D}\Phi_{[\mathbf{A}_{k+1}]}$
- 9: $k \leftarrow k + 1$
- 10: **return** \mathbf{B}_k

random initialization). Indeed, if \mathbf{A}_0 is set to a stationary point of Φ , then $D\Phi_{[\mathbf{A}_0]} = \mathbf{0}$ and, consequently $\tilde{D}\Phi_{[\mathbf{A}_0]} \approx \mathbf{0}$ (given that the approximate gradient is reasonably accurate), which may result in premature and undesirable termination. Clearly, this is not a concern for Algorithm 1 since it assumes convexity of Φ , whereby any stationary point is a global optimum.

The following result follows by adapting the proofs of Theorem 2 and Corollary 2 in [21]. For completeness, we provide a self-contained argument in Appendix D along with a discussion of how this problem fits in the framework of [21].

Theorem 3 (Adaptive convergence rate). *Assume that Φ is L -smooth on \mathcal{D}_M and that $\tilde{\Pi}^A$ is a δ -oracle for Π^A . Then, the iterates $\mathbf{A}_k, \mathbf{B}_k$ in Algorithm 2 with $\tilde{D}\Phi_{[\mathbf{A}_k]}$ given by (6) satisfy*

1. *If Φ is non-convex and $\text{OT}_{(\cdot),\varepsilon}(\mu_0, \mu_1)$ is L' -smooth, then $\min_{1 \leq i \leq k} \|\beta_i^{-1}(\mathbf{B}_i - \mathbf{A}_i)\|_F^2 \leq \frac{96L^2}{k(k+1)(k+2)} \|\mathbf{C}_0 - \mathbf{B}^*\|_F^2 + \frac{24LL'}{k} \left(\|\mathbf{B}^*\|_F^2 + \frac{5M^2 d_0^2 d_1^2}{16} \right) + 8L\delta'$, where \mathbf{B}^* is a global minimizer of Φ , and $\delta' = 32M\delta \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} \|x^{(i)}(y^{(j)})^\top\|_1$.*
2. *If Φ is convex, then $\min_{1 \leq i \leq k} \|\beta_i^{-1}(\mathbf{B}_i - \mathbf{A}_i)\|_F^2 \leq \frac{96L^2}{k(k+1)(k+2)} \|\mathbf{C}_0 - \mathbf{B}^*\|_F^2 + 8L\delta'$.*

We first show that when $\|\beta_k^{-1}(\mathbf{B}_k - \mathbf{A}_k)\|_F$ is small, $D\Phi_{[\mathbf{A}_k]}$ is approximately stationary.

Corollary 2 (Approximate stationarity). *Let $\mathbf{A}_k, \mathbf{B}_k$ be iterates from Algorithm 2 and assume that $\mathbf{B}_k \in \text{int}(\mathcal{D}_M)$. Then, $\|D\Phi_{[\mathbf{A}_k]}\|_F < 32\delta \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} \|x^{(i)}\| \|y^{(j)}\| + \|\beta_k^{-1}(\mathbf{B}_k - \mathbf{A}_k)\|_F$.*

The proof of Corollary 2 follows from the δ -oracle assumption and the fact that when \mathbf{B}_k is an interior point of \mathcal{D}_M , we have $\mathbf{B}_k = \mathbf{A}_k - \beta_k \mathbf{G}_k$. See Appendix A.5 for details. When \mathbf{B}_k is not an interior point of \mathcal{D}_M , the interpretation of $\|\beta_k^{-1}(\mathbf{B}_k - \mathbf{A}_k)\|_F$ is less straightforward. However, as all stationary points of Φ are contained in $\mathcal{D}_{M_{\mu_0, \mu_1}}$, it is expected that Algorithm 2 will converge to an interior point. By analogy with Remark 1, when all iterates are interior points Algorithm 2 yields a bound on the number of iterations required to achieve an approximate stationary point.

The following remark addresses the distinctions between the convex and non-convex settings.

Remark 2 (Adaptivity of Algorithm 2). *As in Theorem 2, the convergence rates are decoupled into a term related to the progress of the optimization procedure and a term related to the oracle error. In the case where Φ is non-convex, the dominant term in the optimization error is $O(1/k)$, which coincides with the best known rates for solving general unconstrained nonlinear programs [21]. On the other hand, when Φ is convex, the rate of convergence improves to $O(1/k^3)$ which essentially matches the best known rates for the norm of the gradient in the unconstrained accelerated gradient method applied to a convex L -smooth function (see Theorem 6 in [34] and Theorem 3.1 in [10]). This adaptivity is beneficial, as Φ may be convex beyond the conditions derived in Theorem 1.*

We now discuss an optimization procedure which does not require convexity of the objective. This accounts for the fact that outside the sufficient conditions of Theorem 1, convexity of Φ is generally unknown. However, the same result shows that Φ is L -smooth with $L = 64 \vee \left(32^2 \varepsilon^{-1} \sqrt{M_4(\mu_0)M_4(\mu_1)} - 64 \right)$ and $\text{OT}_{(\cdot),\varepsilon}$ is L' -smooth with $L' = 32^2 \varepsilon^{-1} \sqrt{M_4(\mu_0)M_4(\mu_1)}$. Hence, we adapt the smooth non-convex optimization routine of [21] to account for our inexact oracle. Notably, their method adapts to the convexity of Φ as described in Theorem 3.

Unlike Algorithm 1, which is initialized at any fixed \mathbf{A}_0 , the starting point in Algorithm 2 should be chosen according to some selection rule that avoids initializing at a stationary point (e.g., ran-

An empirical comparison of Algorithms 1 and 2 in the convex setting is included in Section 4.4. In particular, Algorithm 1 is seen to outperform Algorithm 2 in terms of average runtime despite having the same per iteration complexity when the inexact gradient is computed using Sinkhorn iterations.

Remark 3 (Computational complexity of Algorithms 1 and 2). *As Sinkhorn’s algorithm is known to have a complexity of $O(N_0N_1)$ (cf. e.g. [32]), the gradient approximation (6) can be computed in $O(N_0N_1)$ time. It follows that Algorithms 1 and 2 admit a computational complexity of $O(N_0N_1)$.*

4.4 Numerical Experiments

We conclude this section with some experiments that empirically validate the rates obtained in Theorems 2 and 3 and the computational complexity discussed in Remark 3. All experiments were performed on a desktop computer with 16 GB of RAM and an Intel i5-10600k CPU using the Python programming language. Blown-up figures are included in Appendix F

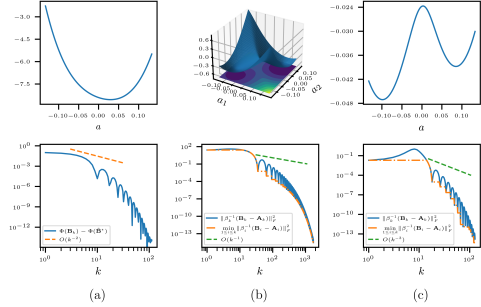


Figure 1: The top row compiles plots of Φ for the different examples. The bottom row consists of plots tracking the progress of the iterates. In (b) and (c), Algorithm 2 is initialized at $\mathbf{C}_0 = (1, 1) \times 10^{-5}$ and $\mathbf{C}_0 = 1 \times 10^{-5}$, respectively.

$\mu_0 = \frac{1}{5}(\delta_{-0.1} + \delta_{-0.2} + \delta_{0.2} + \delta_{-0.3} + \delta_{0.3})$ and $\mu_1 = \frac{1}{5}(\delta_{0.2} + \delta_{-0.3} + \delta_{0.3} + \delta_{-0.4} + \delta_{0.4})$ and $\varepsilon = 0.03$. The stopping condition used in all these example is $\|\mathbf{G}_k\|_F < 5 \times 10^{-8}$ and the approximate gradient (6) is computed using the implementation of Sinkhorn’s algorithm from [19].

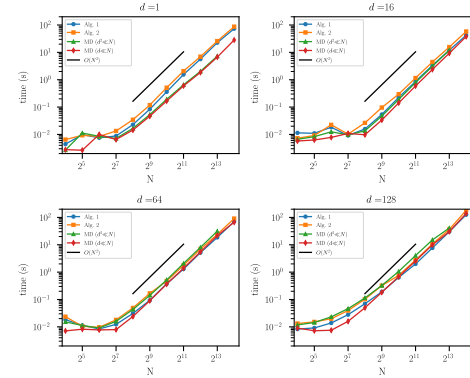


Figure 2: The various plots compile the average runtime of Algorithms 1 and 2, and two versions of the mirror descent algorithm in the convex regime for different combinations of d and N .

The convex case: First, ε is chosen as $1.05 \times 16 \sqrt{M_4(\mu_0)M_4(\mu_1)}$ so as to guarantee convexity of Φ for each instance by Theorem 1. Figure 2 presents the average runtime of both algorithms in this setting with the stopping condition $\|\mathbf{G}_k\|_F < 10^{-6}$. We compare the performance of our

¹The plot shows the approximate gap $\Phi(\mathbf{B}_k) - \Phi(\bar{\mathbf{B}}^*)$, where $\bar{\mathbf{B}}^*$ is the approximate minimizer.

methods with the two implementations of the $O(N^2)$ mirror descent algorithm provided in [32]. The first implementation includes certain algorithmic tweaks when $d^2 \ll N$, whereas the second only requires $d \ll N$ to achieve the quadratic complexity. Our implementation of the mirror descent algorithm is based on the code provided in [32] with some small modifications (e.g., EOT couplings are computed using Sinkhorn’s algorithm from the Python Optimal Transport package [19] and some extraneous logging features are removed). We note that the first version of the mirror descent algorithm encounters “out of memory” errors for $N = 16384$.

The plots show that the four algorithms perform similarly on the considered examples, and empirically validate the $O(N^2)$ computational complexity from Remark 3. To verify that the algorithms all converge to solutions with similar objective values, we evaluate the relative error² between all pairs of algorithms for each d, N . The largest relative error we observe is 6.6×10^{-6} for $d = 1$ and, for the other choices of d , is at most 4.2×10^{-12} . The values obtained are thus in good agreement.

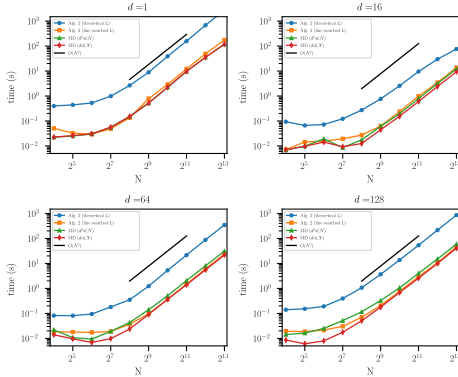


Figure 3: The various plots compile the average runtime of Algorithm 2 with the two methods for choosing L , and two versions of the mirror descent algorithm in the non-convex regime for different combinations of d and N .

algorithm no longer converges. For each d and N , we choose the value of L that attains the fastest convergence, and repeat this procedure for 5 pairs of distributions. For Algorithm 2 with the choice of L that comes from the theoretical bound and the two versions of mirror descent we follow the same methodology as in the convex case. The average runtimes of all methods are reported in Figure 3. The restriction to 5 runs in the line search case is only out of convenience and we note that all algorithms yield similar results if we restrict to 5 runs throughout.

The plots again validate the $O(N^2)$ time complexity for all four approaches. However, we see that choosing L in Algorithm 2 according to the theoretical upper bound may indeed be too conservative, as it results in a $10\times$ slowdown compared to the other methods. By setting L via the line search, on the other hand, Algorithm 2 and mirror descent exhibit similar performance. This suggests that the longer runtime of Algorithm 2 with the theoretical L value can be attributed to this being an overly conservative choice as opposed to a fundamental limitation of this method. Optimization routines that update L at each iteration have been proposed in [3, 26, 37], but require solving an additional EOT problem at each step for our application. As such, these approaches may reduce the number of iterations required for convergence, but at the cost of increasing the per iteration complexity.

5 Conclusion

In this work, we have addressed stability for the EGW problem over Euclidean spaces with quadratic cost. The analysis leveraged variational characterizations of these EGW distances that tie them to EOT with a certain parametrized cost function. The stability analysis was used to study convexity

²Relative error is measured by $\max_{i \in \mathcal{C}(d, N)} |S_\varepsilon^{A1}(\mu_{0,i}, \mu_{1,i}) - S_\varepsilon^{A2}(\mu_{0,i}, \mu_{1,i})| / (S_\varepsilon^{A1}(\mu_{0,i}, \mu_{1,i}) \wedge S_\varepsilon^{A2}(\mu_{0,i}, \mu_{1,i}))$, where $S_\varepsilon^{A1}(\mu_{0,i}, \mu_{1,i})$ and $S_\varepsilon^{A2}(\mu_{0,i}, \mu_{1,i})$ denote the objective values achieved by the first and second algorithm of the pair, and $\mathcal{C}(d, N)$ is the collection of completed runs from a given experiment.

and smoothness properties of this variational problem, which led to two new efficient algorithms for computing the EGW distance. The complexity of these algorithms agrees with the best known complexity of $O(N^2)$ for computing the quadratic EGW distance directly, but unlike previous approaches, we provide, for the first time, non-asymptotic convergence rate guarantees in both the convex and non-convex regimes. This stability analysis also lays the groundwork for solving the EGW problem via smooth optimization methods.

Acknowledgments and Disclosure of Funding

Z. Goldfeld is partially supported by NSF grants CCF-2046018, DMS-2210368, and CCF-2308446, and the IBM Academic Award. K. Kato is partially supported by the NSF grants DMS-1952306, DMS-2014636, and DMS-2210368. G. Rioux is partially supported by the NSERC Postgraduate Fellowship PGSD-567921-2022.

References

- [1] David Alvarez-Melis and Tommi Jaakkola, *Gromov-wasserstein alignment of word embedding spaces*, Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, 2018, pp. 1881–1890.
- [2] Tom M. Apostol, *Mathematical analysis*, 5 ed., Addison-Wesley, 1974.
- [3] Stephen R Becker, Emmanuel J Candès, and Michael C Grant, *Templates for convex cone problems with applications to sparse signal recovery*, Mathematical programming computation **3** (2011), 165–218.
- [4] Garrett Birkhoff, *Extensions of Jentzsch’s theorem*, Transactions of the American Mathematical Society **85** (1957), no. 1, 219–227.
- [5] Andrew J Blumberg, Mathieu Carriere, Michael A Mandell, Raul Rabadan, and Soledad Villar, *MREC: a fast and versatile framework for aligning and matching point clouds with applications to single cell molecular data*, arXiv preprint arXiv:2001.01666 (2020).
- [6] J Frédéric Bonnans and Alexander Shapiro, *Perturbation analysis of optimization problems*, Springer Science & Business Media, 2013.
- [7] Haim Brézis, *Functional analysis, sobolev spaces and partial differential equations*, vol. 2, Springer, 2011.
- [8] Charlotte Bunne, David Alvarez-Melis, Andreas Krause, and Stefanie Jegelka, *Learning generative models across incomparable spaces*, International conference on machine learning, PMLR, 2019, pp. 851–861.
- [9] Guillaume Carlier and Maxime Laborde, *A differential approach to the multi-marginal Schrödinger system*, SIAM Journal on Mathematical Analysis **52** (2020), no. 1, 709–717.
- [10] Shuo Chen, Bin Shi, and Ya-xiang Yuan, *Gradient norm minimization of Nesterov acceleration: $o(1/k^3)$* , arXiv preprint arXiv:2209.08862 (2022).
- [11] Clayton W. Commander, *A survey of the quadratic assignment problem, with applications*, Morehead Electronic Journal of Applicable Mathematics **4** (2005), MATH–2005–01.
- [12] Marco Cuturi, *Sinkhorn distances: lightspeed computation of optimal transport*, Proceedings of the 26th International Conference on Neural Information Processing Systems, 2013, pp. 2292–2300.
- [13] Alexandre d’Aspremont, *Smooth optimization with approximate gradient*, SIAM Journal on Optimization **19** (2008), no. 3, 1171–1183.
- [14] Alexandre d’Aspremont, Laurent El Ghaoui, Michael I. Jordan, and Gert R. G. Lanckriet, *A direct formulation for sparse PCA using semidefinite programming*, SIAM Review **49** (2007), no. 3, 434–448.
- [15] Pinar Demetci, Rebecca Santorella, Björn Sandstede, William Stafford Noble, and Ritambhara Singh, *SCOT: single-cell multi-omics alignment with optimal transport*, Journal of Computational Biology **29** (2022), no. 1, 3–18.

- [16] Olivier Devolder, François Glineur, and Yurii Nesterov, *First-order methods of smooth convex optimization with inexact oracle*, *Mathematical Programming* **146** (2014), 37–75.
- [17] Pavel Dvurechensky, *Gradient method with inexact oracle for composite non-convex optimization*, arXiv preprint arXiv:1703.09180 (2017).
- [18] Pavel Dvurechensky, Alexander Gasnikov, and Alexey Kroshnin, *Computational optimal transport: complexity by accelerated gradient descent is better than by Sinkhorn’s algorithm*, *International conference on machine learning*, PMLR, 2018, pp. 1367–1376.
- [19] Rémi Flamary, Nicolas Courty, Alexandre Gramfort, Mokhtar Z. Alaya, Aurélie Boisbunon, Stanislas Chambon, Laetitia Chapel, Adrien Corenflos, Kilian Fatras, Nemo Fournier, Léo Gautheron, Nathalie T.H. Gayraud, Hicham Janati, Alain Rakotomamonjy, Ievgen Redko, Antoine Rolet, Antony Schutz, Vivien Seguy, Danica J. Sutherland, Romain Tavenard, Alexander Tong, and Titouan Vayer, *POT: Python optimal transport*, *Journal of Machine Learning Research* **22** (2021), no. 78, 1–8.
- [20] Joel Franklin and Jens Lorenz, *On the scaling of multidimensional matrices*, *Linear Algebra and its applications* **114** (1989), 717–735.
- [21] Saeed Ghadimi and Guanghui Lan, *Accelerated gradient methods for nonconvex nonlinear and stochastic programming*, *Mathematical Programming* **156** (2016), no. 1-2, 59–99.
- [22] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal, *Fundamentals of convex analysis*, Springer Science & Business Media, 2004.
- [23] Patrice Koehl, Marc Delarue, and Henri Orland, *Computing the Gromov-Wasserstein distance between two surface meshes using optimal transport*, *Algorithms* **16** (2023), no. 3, 131.
- [24] Facundo Mémoli, *Spectral Gromov-Wasserstein distances for shape matching*, 2009 IEEE 12th International Conference on Computer Vision Workshops, ICCV Workshops, IEEE, 2009, pp. 256–263.
- [25] ———, *Gromov-Wasserstein distances and the metric approach to object matching*, *Found. Comput. Math.* **11** (2011), no. 4, 417–487.
- [26] Yu Nesterov, *Gradient methods for minimizing composite functions*, *Mathematical programming* **140** (2013), no. 1, 125–161.
- [27] Yurii Nesterov, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer Science & Business Media, 2003.
- [28] Marcel Nutz, *Introduction to entropic optimal transport*, *Lecture notes*, Columbia University (2021).
- [29] Gabriel Peyré, Marco Cuturi, and Justin Solomon, *Gromov-Wasserstein averaging of kernel and distance matrices*, *International Conference on Machine Learning*, PMLR, 2016, pp. 2664–2672.
- [30] R Tyrrell Rockafellar, *Convex analysis*, vol. 11, Princeton university press, 1997.
- [31] Hans Samelson, *On the Perron-Frobenius theorem.*, *Michigan Mathematical Journal* **4** (1957), no. 1, 57 – 59.
- [32] Meyer Scetbon, Gabriel Peyré, and Marco Cuturi, *Linear-time Gromov- Wasserstein distances using low rank couplings and costs*, *International Conference on Machine Learning*, PMLR, 2022, pp. 19347–19365.
- [33] Thibault Séjourné, François-Xavier Vialard, and Gabriel Peyré, *The unbalanced Gromov-Wasserstein distance: conic formulation and relaxation*, *Advances in Neural Information Processing Systems* **34** (2021), 8766–8779.
- [34] Bin Shi, Simon S Du, Michael I Jordan, and Weijie J Su, *Understanding the acceleration phenomenon via high-resolution differential equations*, *Mathematical Programming* (2021), 1–70.
- [35] Richard Sinkhorn, *Diagonal equivalence to matrices with prescribed row and column sums*, *The American Mathematical Monthly* **74** (1967), no. 4, 402–405.
- [36] Justin Solomon, Gabriel Peyré, Vladimir G Kim, and Suvrit Sra, *Entropic metric alignment for correspondence problems*, *ACM Transactions on Graphics (ToG)* **35** (2016), no. 4, 1–13.

- [37] Paul Tseng, *On accelerated proximal gradient methods for convex-concave optimization*, <https://www.mit.edu/~dimitrib/PTseng/papers/apgm.pdf>, 2008.
- [38] Titouan Vayer, Rémi Flamary, Romain Tavenard, Laetitia Chapel, and Nicolas Courty, *Sliced Gromov-Wasserstein*, arXiv preprint arXiv:1905.10124 (2020).
- [39] Hongteng Xu, Dixin Luo, and Lawrence Carin, *Scalable Gromov-Wasserstein learning for graph partitioning and matching*, Advances in neural information processing systems **32** (2019).
- [40] Hongteng Xu, Dixin Luo, Hongyuan Zha, and Lawrence Carin Duke, *Gromov-Wasserstein learning for graph matching and node embedding*, International conference on machine learning, PMLR, 2019, pp. 6932–6941.
- [41] Yuguang Yan, Wen Li, Hanrui Wu, Huaqing Min, Mingkui Tan, and Qingyao Wu, *Semi-supervised optimal transport for heterogeneous domain adaptation.*, IJCAI, vol. 7, 2018, pp. 2969–2975.
- [42] Kōsaku Yosida, *Functional analysis*, Springer Science & Business Media, 1995.
- [43] Zhengxin Zhang, Ziv Goldfeld, Youssef Mroueh, and Bharath K. Sriperumbudur, *Gromov-Wasserstein distances: entropic regularization, duality, and sample complexity*, 2022.
- [44] Zhengxin Zhang, Youssef Mroueh, Ziv Goldfeld, and Bharath Sriperumbudur, *Cycle consistent probability divergences across different spaces*, International Conference on Artificial Intelligence and Statistics, PMLR, 2022, pp. 7257–7285.

A Proofs

A.1 Proof of Proposition 1

We first fix some notation. Let $S_i = \text{spt}(\mu_i)$ for $i = 0, 1$ and define the Banach spaces

$$\mathfrak{E} = \left\{ (f_0, f_1) \in \mathcal{C}(S_0) \times \mathcal{C}(S_1) : \int f_0 d\mu_0 = 0 \right\},$$

$$\mathfrak{F} = \left\{ (f_0, f_1) \in \mathcal{C}(S_0) \times \mathcal{C}(S_1) : \int f_0 d\mu_0 = \int f_1 d\mu_1 \right\}.$$

Consider the map $\Upsilon : \mathbb{R}^{d_0 \times d_1} \times \mathfrak{E} \rightarrow \mathcal{C}(S_0) \times \mathcal{C}(S_1)$ given by

$$\Upsilon : (\mathbf{A}, \varphi_0, \varphi_1) \mapsto \left(\int e^{\frac{\varphi_0(\cdot) + \varphi_1(y) - c_{\mathbf{A}}(\cdot, y)}{\varepsilon}} d\mu_1(y) - 1, \int e^{\frac{\varphi_0(x) + \varphi_1(\cdot) - c_{\mathbf{A}}(x, \cdot)}{\varepsilon}} d\mu_0(x) - 1 \right).$$

For fixed $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$, the solution to the equation $\Upsilon(\mathbf{A}, \cdot, \cdot) = 0$ is the unique pair of EOT potentials $(\varphi_0^{\mathbf{A}}, \varphi_1^{\mathbf{A}})$ for μ_0, μ_1 with the cost $c_{\mathbf{A}}$ satisfying the normalization from \mathfrak{E} . Observe that, by compactness of S_0 and S_1 , the potentials are bounded.

The following lemmas verify the conditions to apply the implicit mapping theorem to Υ in order to obtain the Fréchet derivative of the map $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1} \mapsto (\varphi_0^{\mathbf{A}}, \varphi_1^{\mathbf{A}})$. Given that $\text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1) = \int \varphi_0^{\mathbf{A}} d\mu_0 + \int \varphi_1^{\mathbf{A}} d\mu_1$, the derivative of the map $\mathbf{A} \mapsto \text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1)$ and that of Φ itself will readily follow.

Lemma 1. *The map Υ is smooth with first derivative at $(\mathbf{A}, \varphi_0, \varphi_1) \in \mathbb{R}^{d_0 \times d_1} \times \mathfrak{E}$ given by,*

$$D\Upsilon_{[\mathbf{A}, \varphi_0, \varphi_1]}(\mathbf{B}, h_0, h_1) = \varepsilon^{-1} \left(\int (h_0(\cdot) + h_1(y) + 32(\cdot)^\top \mathbf{B}y) e^{\frac{\varphi_0(\cdot) + \varphi_1(y) - c_{\mathbf{A}}(\cdot, y)}{\varepsilon}} d\mu_1(y), \right. \\ \left. \int (h_0(x) + h_1(\cdot) + 32x^\top \mathbf{B}(\cdot)) e^{\frac{\varphi_0(x) + \varphi_1(\cdot) - c_{\mathbf{A}}(x, \cdot)}{\varepsilon}} d\mu_0(x) \right),$$

where $(\mathbf{B}, h_0, h_1) \in \mathbb{R}^{d_0 \times d_1} \times \mathfrak{E}$.

The proof of this result is straightforward, but included in Appendix E.1 for completeness. Now, define $\xi_{\mathbf{A}} := \varepsilon D\Upsilon_{[\mathbf{A}, \varphi_0^{\mathbf{A}}, \varphi_1^{\mathbf{A}}]}(0, \cdot, \cdot)$ and let $\pi_{\mathbf{A}}$ be the EOT coupling for $\text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1)$. Note that for any $(h_0, h_1) \in \mathfrak{E}$, we have $\xi_{\mathbf{A}}(h_0, h_1) \in \mathfrak{F}$, which follows by recalling that $\frac{d\pi_{\mathbf{A}}}{d\mu_0 \otimes \mu_1}(x, y) = e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x, y)}{\varepsilon}}$ and observing

$$\int (\xi_{\mathbf{A}}(h_0, h_1))_1 d\mu_0 = \int h_0 d\mu_0 + \int h_1 d\pi_{\mathbf{A}} = \int h_0 d\mu_0 + \int h_1 d\mu_1$$

$$\int (\xi_{\mathbf{A}}(h_0, h_1))_2 d\mu_1 = \int h_0 d\pi_{\mathbf{A}} + \int h_1 d\mu_1 = \int h_0 d\mu_0 + \int h_1 d\mu_1.$$

We next prove that $\xi_{\mathbf{A}}$ is an isomorphism between \mathfrak{E} and \mathfrak{F} by following the proof of Proposition 3.1 in [9].

Lemma 2. *The map $\xi_{\mathbf{A}}$ is an isomorphism between \mathfrak{E} and \mathfrak{F} .*

Proof. Observe that $\xi_{\mathbf{A}}$ extends naturally to a map on $L^2(\mu_0) \times L^2(\mu_1)$ and admits the representation

$$\xi_{\mathbf{A}} : (f_0, f_1) \in L^2(\mu_0) \times L^2(\mu_1) \mapsto (f_0, f_1) + \mathcal{L}(f_0, f_1) \in L^2(\mu_0) \times L^2(\mu_1),$$

where

$$\mathcal{L}(f_0, f_1) = \left(\int f_1(y) e^{\frac{\varphi_0^{\mathbf{A}}(\cdot) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(\cdot, y)}{\varepsilon}} d\mu_1(y), \int f_0(x) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(\cdot) - c_{\mathbf{A}}(x, \cdot)}{\varepsilon}} d\mu_0(x) \right).$$

Lemma 11 in Appendix E.2 demonstrates that \mathcal{L} is a compact linear self-map of $L^2(\mu_0) \times L^2(\mu_1)$.

We first show that $\xi_{\mathbf{A}}$ is injective on $E := \{(f_0, f_1) \in L^2(\mu_0) \times L^2(\mu_1) : \int f_0 d\mu_0 = 0\}$. Suppose that (\bar{f}_0, \bar{f}_1) satisfies $\xi_{\mathbf{A}}(\bar{f}_0, \bar{f}_1) = 0$. Multiplying $(\xi_{\mathbf{A}}(\bar{f}_0, \bar{f}_1))_1$ by \bar{f}_0 and $(\xi_{\mathbf{A}}(\bar{f}_0, \bar{f}_1))_2$ by \bar{f}_1 , we have

$$\begin{aligned} \int (\bar{f}_0^2(\cdot) + \bar{f}_0(\cdot)\bar{f}_1(y)) e^{\frac{\varphi_0^{\mathbf{A}}(\cdot) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(\cdot, y)}{\varepsilon}} d\mu_1(y) &= 0, \\ \int (\bar{f}_0(x)f_1(\cdot) + \bar{f}_1^2(\cdot)) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(\cdot) - c_{\mathbf{A}}(x, \cdot)}{\varepsilon}} d\mu_0(x) &= 0, \end{aligned}$$

and summing these equations gives $\int (\bar{f}_0 + \bar{f}_1)^2 d\pi_{\mathbf{A}} = 0$. As $\pi_{\mathbf{A}}$ is equivalent to $\mu_0 \otimes \mu_1$, we have $\bar{f}_0 + \bar{f}_1 = 0$ $\mu_0 \otimes \mu_1$ -a.e., which further implies that $(\bar{f}_0, \bar{f}_1) = (a, -a)$ $\mu_0 \otimes \mu_1$ -a.e. for some $a \in \mathbb{R}$. Consequently, $\ker(\xi_{\mathbf{A}})$ is 1-dimensional and $\xi_{\mathbf{A}}$ is injective on E .

Next, we show that $\xi_{\mathbf{A}}$ is onto $F := \{(f_0, f_1) \in L^2(\mu_0) \times L^2(\mu_1) : \int f_0 d\mu_0 = \int f_1 d\mu_1\}$. As in the lead-up to this lemma, $\xi_{\mathbf{A}}(E) \subset F$. By the Fredholm alternative (cf. Theorem 6.6 in [7]), $(\text{Id} + \mathcal{L})(L^2(\mu_0) \times L^2(\mu_1))$ has codimension 1 and, as F has codimension 1, we must have $\xi_{\mathbf{A}}(E) = F$.

As such, for any $(g_0, g_1) \in \mathfrak{F} \subset F$, there exists $(h_0, h_1) \in E$ for which

$$\xi_{\mathbf{A}}(h_0, h_1) = (h_0, h_1) + \mathcal{L}(h_0, h_1) = (g_0, g_1).$$

As $\mathcal{L}(h_0, h_1) \in \mathcal{C}(S_0) \times \mathcal{C}(S_1)$, $(h_0, h_1) = (g_0, g_1) - \mathcal{L}(h_0, h_1) \in \mathcal{C}(S_0) \times \mathcal{C}(S_1)$ with $\int h_0 d\mu_0 = 0$, and thus $(h_0, h_1) \in \mathfrak{E}$. This implies that $\xi_{\mathbf{A}}(\mathfrak{E}) \supset \mathfrak{F}$ and from before we have $\xi_{\mathbf{A}}(\mathfrak{E}) \subset \mathfrak{F}$, yielding $\xi_{\mathbf{A}}(\mathfrak{E}) = \mathfrak{F}$. We have shown that $\xi_{\mathbf{A}} : \mathfrak{E} \rightarrow \mathfrak{F}$ is a continuous linear bijection and hence an isomorphism by the open mapping theorem (cf. Corollary 2.7 in [7]). \square

We now apply the implicit mapping theorem to obtain the Fréchet derivative of $(\varphi_0^{(\cdot)}, \varphi_1^{(\cdot)})$.

Lemma 3. *The map $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1} \mapsto (\varphi_0^{\mathbf{A}}, \varphi_1^{\mathbf{A}}) \in \mathfrak{E}$ is smooth with Fréchet derivative*

$$D\left(\varphi_0^{(\cdot)}, \varphi_1^{(\cdot)}\right)_{[\mathbf{A}]}(\mathbf{B}) = -\left(h_0^{\mathbf{A}, \mathbf{B}}, h_1^{\mathbf{A}, \mathbf{B}}\right),$$

where $(h_0^{\mathbf{A}, \mathbf{B}}, h_1^{\mathbf{A}, \mathbf{B}}) \in \mathfrak{E}$ satisfies

$$\begin{aligned} \int \left(h_0^{\mathbf{A}, \mathbf{B}}(x) + h_1^{\mathbf{A}, \mathbf{B}}(y) - 32x^{\top} \mathbf{B}y\right) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x, y)}{\varepsilon}} d\mu_1(y) &= 0, \quad \forall x \in \text{spt}(\mu_0), \\ \int \left(h_0^{\mathbf{A}, \mathbf{B}}(x) + h_1^{\mathbf{A}, \mathbf{B}}(y) - 32x^{\top} \mathbf{B}y\right) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x, y)}{\varepsilon}} d\mu_0(x) &= 0, \quad \forall y \in \text{spt}(\mu_1), \end{aligned} \quad (7)$$

with $(\varphi_0^{\mathbf{A}}, \varphi_1^{\mathbf{A}})$ being any pair of EOT potentials for (μ_0, μ_1) with the cost $c_{\mathbf{A}}$.

Proof. Fix $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$ with corresponding EOT potentials $(\varphi_0^{\mathbf{A}}, \varphi_1^{\mathbf{A}})$. For notational convenience, define the shorthands $D_1 \Upsilon_{\mathbf{A}} = D\Upsilon_{[\mathbf{A}, \varphi_0^{\mathbf{A}}, \varphi_1^{\mathbf{A}}]}(\cdot, 0, 0)$ and $D_2 \Upsilon_{\mathbf{A}} = D\Upsilon_{[\mathbf{A}, \varphi_0^{\mathbf{A}}, \varphi_1^{\mathbf{A}}]}(0, \cdot, \cdot)$ (cf. Lemma 1). By Lemma 2, $D_2 \Upsilon_{\mathbf{A}}$ is an isomorphism and we may invoke the implicit mapping theorem (cf. Theorem 5.14 in [6]). This implies that there exists an open neighborhood $U \subset \mathbb{R}^{d_0 \times d_1}$ of \mathbf{A} and a smooth map $g : U \rightarrow \mathfrak{E}$ for which $\Upsilon(\mathbf{B}, g(\mathbf{B})) = 0$ for every $\mathbf{B} \in U$ and

$$Dg_{[\mathbf{A}]}(\mathbf{B}) = -(D_2 \Upsilon_{\mathbf{A}})^{-1} (D_1 \Upsilon_{\mathbf{A}}(\mathbf{B})),$$

i.e., $-Dg_{[\mathbf{A}]}(\mathbf{B})$ solves (7). By construction, $g(\mathbf{B}) = (\varphi_0^{\mathbf{B}}, \varphi_1^{\mathbf{B}})$ and by repeating this process at any $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$, we extend the differentiability of the potentials to the entire space $\mathbb{R}^{d_0 \times d_1}$. \square

Given the dual form of the EOT cost, Lemma 3 suffices to prove Proposition 1.

Proof of Proposition 1. As $\text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1) = \int \varphi_0^{\mathbf{A}} d\mu_0 + \int \varphi_1^{\mathbf{A}} d\mu_1$, Lemma 3 implies that $\text{OT}_{(\cdot), \varepsilon}(\mu_0, \mu_1)$ is smooth with first derivative at $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$ given by

$$D(\text{OT}_{(\cdot), \varepsilon}(\mu_0, \mu_1))_{[\mathbf{A}]}(\mathbf{B}) = -\int h_0^{\mathbf{A}, \mathbf{B}} d\mu_0 - \int h_1^{\mathbf{A}, \mathbf{B}} d\mu_1,$$

where $\mathbf{B} \in \mathbb{R}^{d_0 \times d_1}$. Integrating the first equation in (7) w.r.t. μ_0 while using $\frac{d\pi_{\mathbf{A}}}{d\mu_0 \otimes \mu_1}(x, y) = e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x, y)}{\varepsilon}}$, yields

$$\int \left(h_0^{\mathbf{A}, \mathbf{B}}(x) + h_1^{\mathbf{A}, \mathbf{B}}(y) \right) d\pi_{\mathbf{A}}(x, y) = \int h_0^{\mathbf{A}, \mathbf{B}} d\mu_0 + \int h_1^{\mathbf{A}, \mathbf{B}} d\mu_1 = 32 \int x^\top \mathbf{B} y d\pi_{\mathbf{A}}(x, y), \quad (8)$$

whence

$$D(\text{OT}_{(\cdot, \varepsilon)}(\mu_0, \mu_1))_{[\mathbf{A}]}(\mathbf{B}) = -32 \int x^\top \mathbf{B} y d\pi_{\mathbf{A}}(x, y).$$

As $\|\mathbf{A}\|_F^2 = \text{tr}(\mathbf{A}^\top \mathbf{A})$, we have $D(32\|\cdot\|_F^2)_{[\mathbf{A}]}(\mathbf{B}) = 64\text{tr}(\mathbf{A}^\top \mathbf{B})$, which together with the display above yields

$$D\Phi_{[\mathbf{A}]}(\mathbf{B}) = 64 \text{tr}(\mathbf{A}^\top \mathbf{B}) - 32 \int x^\top \mathbf{B} y d\pi_{\mathbf{A}}(x, y),$$

as desired.

For the second-order derivative, recall from Section 2.1 that $\frac{d\pi_{\mathbf{A}}}{d\mu_0 \otimes \mu_1}(x, y) = e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x, y)}{\varepsilon}}$. As in the proof of Lemma 1, as the map

$$\mathbf{A} \in \mathbb{R}^{d_0 \times d_1} \mapsto ((x, y) \in S_0 \times S_1 \mapsto \varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x, y)) \in \mathcal{C}(S_0 \times S_1)$$

is differentiable at $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$ with derivative

$$\mathbf{C} \in \mathbb{R}^{d_0 \times d_1} \mapsto ((x, y) \in S_0 \times S_1 \mapsto (h_0^{\mathbf{A}, \mathbf{C}}(x) + h_1^{\mathbf{A}, \mathbf{C}}(y) - 32x^\top \mathbf{C} y)) \in \mathcal{C}(S_0 \times S_1),$$

the expansion

$$\frac{d\pi_{\mathbf{A}+\mathbf{C}}}{d\mu_0 \otimes \mu_1}(x, y) - \frac{d\pi_{\mathbf{A}}}{d\mu_0 \otimes \mu_1}(x, y) = -\varepsilon^{-1} z_{\mathbf{A}, \mathbf{C}}(x, y) \frac{d\pi_{\mathbf{A}}}{d\mu_0 \otimes \mu_1}(x, y) + R_{\mathbf{C}}(x, y),$$

holds uniformly over $(x, y) \in S_0 \times S_1$, where $R_{\mathbf{C}}(x, y) = o(\mathbf{C})$ as $\|\mathbf{C}\|_F \rightarrow 0$ and $z_{\mathbf{A}, \mathbf{C}}(x, y) = h_0^{\mathbf{A}, \mathbf{C}}(x) + h_1^{\mathbf{A}, \mathbf{C}}(y) - 32x^\top \mathbf{C} y$. Thus,

$$\begin{aligned} & \sup_{\|\mathbf{B}\|_F=1} \frac{\left| \int x^\top \mathbf{B} y d\pi_{\mathbf{A}+\mathbf{C}}(x, y) - \int x^\top \mathbf{B} y d\pi_{\mathbf{A}}(x, y) - \varepsilon^{-1} \int x^\top \mathbf{B} y z_{\mathbf{A}, \mathbf{C}}(x, y) d\pi_{\mathbf{A}}(x, y) \right|}{\|\mathbf{C}\|_F} \\ &= \sup_{\|\mathbf{B}\|_F=1} \left| \int x^\top \mathbf{B} y \|\mathbf{C}\|_F^{-1} R_{\mathbf{C}}(x, y) d\mu_0 \otimes \mu_1(x, y) \right| \\ &\leq \sup_{(x, y) \in S_1 \times S_2} \|x\| \|y\| \int \|\mathbf{C}\|_F^{-1} |R_{\mathbf{C}}(x, y)| d\mu_0 \otimes \mu_1(x, y). \end{aligned}$$

As $R_{\mathbf{C}}(x, y) = o(\mathbf{C})$, this final term converges to 0 as $\|\mathbf{C}\|_F \rightarrow 0$, so

$$D^2(\text{OT}_{(\cdot, \varepsilon)}(\mu_0, \mu_1))_{[\mathbf{A}]}(\mathbf{B}, \mathbf{C}) = 32\varepsilon^{-1} \int x^\top \mathbf{B} y \left(h_0^{\mathbf{A}, \mathbf{C}}(x) + h_1^{\mathbf{A}, \mathbf{C}}(y) - 32x^\top \mathbf{C} y \right) d\pi_{\mathbf{A}}(x, y).$$

As $D(32\|\cdot\|_F^2)_{[\mathbf{A}]}(\mathbf{B}) = 64\text{tr}(\mathbf{A}^\top \mathbf{B})$, $D^2(32\|\cdot\|_F^2)_{[\mathbf{A}]}(\mathbf{B}, \mathbf{C}) = 64\text{tr}(\mathbf{C}^\top \mathbf{B})$. Altogether,

$$D^2\Phi_{[\mathbf{A}]}(\mathbf{B}, \mathbf{C}) = 64 \text{tr}(\mathbf{B}^\top \mathbf{C}) + 32\varepsilon^{-1} \int x^\top \mathbf{B} y \left(h_0^{\mathbf{A}, \mathbf{C}}(x) + h_1^{\mathbf{A}, \mathbf{C}}(y) - 32x^\top \mathbf{C} y \right) d\pi_{\mathbf{A}}(x, y).$$

Coercivity of Φ follows from the fact that

$$\begin{aligned} \text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1) &= \inf_{\pi \in \Pi(\mu_0, \mu_1)} \left\{ \int -4\|x\|^2 \|y\|^2 - 32x^\top \mathbf{A} y d\pi(x, y) + \varepsilon \text{D}_{\text{KL}}(\pi \| \mu_0 \otimes \mu_1) \right\}, \\ &\geq \inf_{\pi \in \Pi(\mu_0, \mu_1)} \left\{ \int -4\|x\|^2 \|y\|^2 - 32\|\mathbf{A}\|_F \|x\| \|y\| d\pi(x, y) \right\} \\ &\geq -4\sqrt{M_4(\mu_0)M_4(\mu_1)} - 32\|\mathbf{A}\|_F \sqrt{M_2(\mu_0)M_2(\mu_1)}, \end{aligned}$$

such that $\Phi(\mathbf{A}) = 32\|\mathbf{A}\|_F^2 + \text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1) \rightarrow +\infty$ as $\|\mathbf{A}\|_F \rightarrow \infty$. \square

A.2 Proof of Corollary 1

Item (i). The expression for the stationary points follows immediately from Proposition 1. To see that all stationary points are elements of $\mathcal{D}_{M_{\mu_0, \mu_1}}$, observe that if \mathbf{A} is a stationary point, then

$$|\mathbf{A}_{ij}| = \frac{1}{2} \left| \int x_i y_j d\pi_{\mathbf{A}}(x, y) \right| \leq \frac{1}{2} \int |x_i y_j| d\pi_{\mathbf{A}}(x, y) \leq \frac{1}{2} \sqrt{M_2(\mu_0)M_2(\mu_1)}.$$

Item (ii). As discussed in Section 2.2, if π_* is optimal for S_ε then $\frac{1}{2} \int xy^\top d\pi_*(x, y)$ minimizes Φ . On the other hand, if \mathbf{A} minimizes Φ , then we have $\mathbf{A} = \frac{1}{2} \int xy^\top d\pi_{\mathbf{A}}$ and hence

$$\begin{aligned} S_{2, \varepsilon}(\mu_0, \mu_1) &= 8 \left\| \int xy^\top d\pi_{\mathbf{A}}(x, y) \right\|_F^2 - 4 \int \|x\|^2 \|y\|^2 d\pi_{\mathbf{A}}(x, y) \\ &\quad - 32 \left\langle \frac{1}{2} \int xy^\top d\pi_{\mathbf{A}}, \int xy^\top d\pi_{\mathbf{A}} \right\rangle_F + \varepsilon \text{D}_{\text{KL}}(\pi_{\mathbf{A}} \| \mu_0 \otimes \mu_1) \\ &= -4 \int \|x\|^2 \|y\|^2 d\pi_{\mathbf{A}(x, y)} - 8 \left\| \int xy^\top d\pi_{\mathbf{A}}(x, y) \right\|_F^2 + \varepsilon \text{D}_{\text{KL}}(\pi_{\mathbf{A}} \| \mu_0 \otimes \mu_1). \end{aligned}$$

By (4),

$$\begin{aligned} S_\varepsilon(\mu_0, \mu_1) &= S_\varepsilon(\mu_0, \mu_1) + S_{2, \varepsilon}(\mu_0, \mu_1) \\ &= \int \left| \|x - x'\|^2 - \|y - y'\|^2 \right|^2 + 2\|x - x'\|^2 \|y - y'\|^2 d\pi_{\mathbf{A}} \otimes \pi_{\mathbf{A}}(x, y, x', y') \\ &\quad - 4 \int \|x\|^2 \|y\|^2 d\mu_0 \otimes \mu_1(x, y) - 4 \int \|x\|^2 \|y\|^2 d\pi_{\mathbf{A}}(x, y) \\ &\quad - 8 \left\| \int xy^\top d\pi_{\mathbf{A}}(x, y) \right\|_F^2 + \varepsilon \text{D}_{\text{KL}}(\pi_{\mathbf{A}} \| \mu_0 \otimes \mu_1). \end{aligned} \tag{9}$$

As $\|x - x'\|^2 \|y - y'\|^2 = (\|x\|^2 - 2x^\top x' + \|x'\|^2) (\|y\|^2 - 2y^\top y' + \|y'\|^2)$, we have

$$\begin{aligned} &\int \|x - x'\|^2 \|y - y'\|^2 d\pi_{\mathbf{A}} \otimes \pi_{\mathbf{A}}(x, y, x', y') \\ &= 2 \int \|x\|^2 \|y\|^2 d\mu_0 \otimes \mu_1(x, y) + 2 \int \|x\|^2 \|y\|^2 d\pi_{\mathbf{A}}(x, y) \\ &\quad + 4 \int x^\top x' y^\top y' d\pi_{\mathbf{A}} \otimes \pi_{\mathbf{A}}(x, y, x', y'), \end{aligned}$$

which, together with (9) yields

$$S_\varepsilon(\mu_0, \mu_1) = \int \left| \|x - x'\|^2 - \|y - y'\|^2 \right|^2 d\pi_{\mathbf{A}} \otimes \pi_{\mathbf{A}}(x, y, x', y') + \varepsilon \text{D}_{\text{KL}}(\pi_{\mathbf{A}} \| \mu_0 \otimes \mu_1),$$

proving optimality of $\pi_{\mathbf{A}}$.

Item (iii). Suppose S_ε admits a unique optimal coupling. If two matrices \mathbf{A} and \mathbf{B} minimize Φ , then $\pi_{\mathbf{A}} = \pi_{\mathbf{B}}$ by uniqueness, so $\mathbf{A} = \frac{1}{2} \int xy^\top d\pi_{\mathbf{A}}(x, y) = \frac{1}{2} \int xy^\top d\pi_{\mathbf{B}}(x, y) = \mathbf{B}$. Conversely, suppose Φ admits a unique minimizer \mathbf{A}^* . If π is optimal for S_ε , then π solves the EOT problem $\text{OT}_{\mathbf{A}^*, \varepsilon}(\mu_0, \mu_1)$, so $\pi = \pi_{\mathbf{A}^*}$. \square

A.3 Proof of Theorem 1

The proof of Theorem 1 depends on the following lemma. The variance bound is seen to be sharp up to constants in Appendix B.

Lemma 4 (Hessian eigenvalue bounds). *The following hold for any $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$:*

(i) *The minimal eigenvalue of $D^2\Phi_{[\mathbf{A}]}$, $\lambda_{\min}(D^2\Phi_{[\mathbf{A}]})$, admits the lower bound $64 - 32^2\varepsilon^{-1} \sup_{\|C\|_F=1} \text{Var}_{\pi_{\mathbf{A}}}(X^\top CY)$, where $\sup_{\|C\|_F=1} \text{Var}_{\pi_{\mathbf{A}}}(X^\top CY) \leq \sqrt{M_4(\mu_0)M_4(\mu_1)}$.*

(ii) The maximal eigenvalue of $D^2\Phi_{[\mathbf{A}]}$ satisfies $\lambda_{\max}(D^2\Phi_{[\mathbf{A}]}) \leq 64$.

Proof. We first prove Item (i). The minimal eigenvalue of $D^2\Phi_{[\mathbf{A}]}$ is given in variational form as

$$\begin{aligned} & \inf_{\|\mathbf{C}\|_F=1} D^2\Phi_{[\mathbf{A}]}(\mathbf{C}, \mathbf{C}) \\ &= \inf_{\|\mathbf{C}\|_F=1} \left\{ 64\|\mathbf{C}\|_F^2 + 32\varepsilon^{-1} \int x^\top \mathbf{C} y \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) - 32x^\top \mathbf{C} y \right) d\pi_{\mathbf{A}}(x, y) \right\} \\ &\geq 64 + 32\varepsilon^{-1} \inf_{\|\mathbf{C}\|_F=1} \left\{ \int x^\top \mathbf{C} y \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) - 32x^\top \mathbf{C} y \right) d\pi_{\mathbf{A}}(x, y) \right\}, \end{aligned}$$

using the formula for $D^2\Phi_{[\mathbf{A}]}$ from Proposition 1. Recall that $(h_0^{\mathbf{A},\mathbf{C}}, h_1^{\mathbf{A},\mathbf{C}})$ satisfy

$$\begin{aligned} & \int \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) - 32x^\top \mathbf{C} y \right) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x,y)}{\varepsilon}} d\mu_1(y) = 0, \quad \forall x \in \text{spt}(\mu_0), \\ & \int \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) - 32x^\top \mathbf{C} y \right) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x,y)}{\varepsilon}} d\mu_0(x) = 0, \quad \forall y \in \text{spt}(\mu_1), \end{aligned}$$

such that, multiplying the top equation by $h_0^{\mathbf{A},\mathbf{C}}$ and integrating w.r.t. μ_0 and performing the same operations on the lower equation with $h_1^{\mathbf{A},\mathbf{C}}$ and μ_1 respectively,

$$\begin{aligned} & \int \left[\left(h_0^{\mathbf{A},\mathbf{C}}(x) \right)^2 + h_1^{\mathbf{A},\mathbf{C}}(y) h_0^{\mathbf{A},\mathbf{C}}(x) - 32x^\top \mathbf{C} y h_0^{\mathbf{A},\mathbf{C}}(x) \right] d\pi_{\mathbf{A}}(x, y) = 0, \\ & \int \left[h_0^{\mathbf{A},\mathbf{C}}(x) h_1^{\mathbf{A},\mathbf{C}}(y) + \left(h_1^{\mathbf{A},\mathbf{C}}(y) \right)^2 - 32x^\top \mathbf{C} y h_1^{\mathbf{A},\mathbf{C}}(y) \right] d\pi_{\mathbf{A}}(x, y) = 0. \end{aligned}$$

Summing these equations gives

$$32 \int x^\top \mathbf{C} y \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) \right) d\pi_{\mathbf{A}}(x, y) = \int \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) \right)^2 d\pi_{\mathbf{A}}(x, y),$$

such that

$$\begin{aligned} & 32 \int x^\top \mathbf{C} y \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) - 32x^\top \mathbf{C} y \right) d\pi_{\mathbf{A}}(x, y) \\ &= \int \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) \right)^2 d\pi_{\mathbf{A}}(x, y) - 32^2 \int (x^\top \mathbf{C} y)^2 d\pi_{\mathbf{A}}(x, y), \end{aligned}$$

which proves the first part of Item (i). It remains to show that

$$\int \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) \right)^2 d\pi_{\mathbf{A}}(x, y) - 32^2 \int (x^\top \mathbf{C} y)^2 d\pi_{\mathbf{A}}(x, y) \geq -32^2 \text{Var}_{\pi_{\mathbf{A}}}[X^\top \mathbf{C} Y].$$

By Jensen's inequality, we have

$$\begin{aligned} & \int \left(h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) \right)^2 d\pi_{\mathbf{A}}(x, y) \geq \left(\int h_0^{\mathbf{A},\mathbf{C}}(x) + h_1^{\mathbf{A},\mathbf{C}}(y) d\pi_{\mathbf{A}}(x, y) \right)^2 \\ &= 32^2 \left(\int x^\top \mathbf{C} y d\pi_{\mathbf{A}}(x, y) \right)^2, \end{aligned}$$

where the equality follows from (8), proving the desired inequality.

To prove the uniform bound on the variance in Item (i), observe that

$$\begin{aligned} \sup_{\|\mathbf{C}\|_F=1} \text{Var}_{\pi_{\mathbf{A}}}[X^\top \mathbf{C} Y] &\leq \sup_{\|\mathbf{C}\|_F=1} \mathbb{E}_{\pi_{\mathbf{A}}}[(X^\top \mathbf{C} Y)^2] \\ &\leq \sup_{\|\mathbf{C}\|_F=1} \|\mathbf{C}\|_F^2 \int \|x\|^2 \|y\|^2 d\pi_{\mathbf{A}}(x, y), \\ &\leq \sqrt{M_4(\mu_0)M_4(\mu_1)} \end{aligned}$$

where the final two inequalities follow from the Cauchy-Schwarz inequality.

We now prove the upper bound on the maximum eigenvalue of $D^2\Phi_{[\mathbf{A}]}$ from Item (ii) again using its variational characterization,

$$\lambda_{\max}(D^2\Phi_{[\mathbf{A}]}) = \sup_{\|\mathbf{C}\|_F=1} D^2\Phi_{[\mathbf{A}]}(\mathbf{C}, \mathbf{C}) = 64 + \lambda_{\max}(D^2\text{OT}_{(\cdot),\varepsilon}(\mu_0, \mu_1)_{[\mathbf{A}]}) .$$

Observe that $\text{OT}_{\mathbf{A},\varepsilon}(\mu_0, \mu_1) = \inf_{\pi \in \Pi(\mu_0, \mu_1)} g(\mathbf{A}, \pi, \mu_0, \mu_1, \varepsilon)$, where g depends on \mathbf{A} only through the term $32\text{tr}(\mathbf{A}^\top \int xy^\top d\pi(x, y))$ which is linear in \mathbf{A} . It follows from, e.g., Proposition 2.1.2 in [22] that $\text{OT}_{(\cdot),\varepsilon}(\mu_0, \mu_1)$ is concave. As such, $\lambda_{\max}(D^2\text{OT}_{(\cdot),\varepsilon}(\mu_0, \mu_1)_{[\mathbf{A}]}) \leq 0$, so $\lambda_{\max}(D^2\Phi_{[\mathbf{A}]}) \leq 64$. \square

Proof of Theorem 1. We first discuss the convexity properties of Φ . By Lemma 4, $\lambda_{\min}(D^2\Phi_{[\mathbf{A}]} + \frac{\rho}{2}\|\mathbf{A}\|_F^2) \geq 64 - 32^2\varepsilon^{-1}\sqrt{M_4(\mu_0)M_4(\mu_1)} + \rho$ for any $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$ and $\rho \geq 0$. When this lower bound is non-negative, Φ is ρ -weakly convex on $\mathbb{R}^{d_0 \times d_1}$ by definition. It follows that Φ is always ρ -weakly convex for $\rho = 32^2\varepsilon^{-1}\sqrt{M_4(\mu_0)M_4(\mu_1)} - 64$. Moreover, if $\sqrt{M_4(\mu_0)M_4(\mu_1)} < \frac{\varepsilon}{16}$, then $\lambda_{\min}(D^2\Phi_{[\mathbf{A}]}) > 0$ such that Φ is strictly convex.

L -smoothness of Φ follows from the mean value inequality (see e.g. Example 2 [2, p.356])

$$\begin{aligned} \|D\Phi_{[\mathbf{A}]} - D\Phi_{[\mathbf{B}]}\|_F &\leq \sup_{\mathbf{C} \in [\mathbf{A}, \mathbf{B}]} \sup_{\|\mathbf{E}\|_F=1} |D^2\Phi_{[\mathbf{C}]}(\mathbf{A} - \mathbf{B}, \mathbf{E})|, \\ &\leq \sup_{\mathbf{C} \in [\mathbf{A}, \mathbf{B}]} (|\lambda_{\min}(D^2\Phi_{[\mathbf{C}]})| \vee |\lambda_{\max}(D^2\Phi_{[\mathbf{C}]})|) \|\mathbf{A} - \mathbf{B}\|_F, \end{aligned}$$

for any $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d_0 \times d_1}$, where $[\mathbf{A}, \mathbf{B}]$ denotes the line segment connecting \mathbf{A} and \mathbf{B} . The claimed result then follows by noting that, for any $\mathbf{A}, \mathbf{B} \in \mathcal{D}_M$, $[\mathbf{A}, \mathbf{B}] \subset \mathcal{D}_M$ by convexity and the supremum over \mathcal{D}_M is achieved by compactness and the fact that the objective is continuous. Indeed, the maps $\lambda_{\max}(\cdot)$, $\lambda_{\min}(\cdot)$ are continuous on the space of symmetric matrices, and $D^2\Phi_{[\cdot]}$ is continuous as Φ is smooth. \square

A.4 Proof of Theorem 2

In this section, we show that Theorem 2.2 in [13] on the convergence rate of Algorithm 1 is applicable in our setting. We particularize their result to a fixed prox-function $d = \frac{1}{2}\|\cdot\|_F^2$ which is smooth and 1-strongly convex.

First, we justify the expressions for the iterates $\mathbf{B}_k, \mathbf{C}_k$ in Algorithm 1, which are defined in [13] as the proximal operators

$$\begin{aligned} \mathbf{B}_k &= \operatorname{argmin}_{\mathbf{V} \in \mathcal{D}_M} \left\{ \text{tr}(\mathbf{G}_k^\top \mathbf{V}) + \frac{L}{2} \|\mathbf{V} - \mathbf{A}_k\|_F^2 \right\}, \\ \mathbf{C}_k &= \operatorname{argmin}_{\mathbf{V} \in \mathcal{D}_M} \left\{ \text{tr}(\mathbf{W}_k^\top \mathbf{V}) + \frac{L}{2} \|\mathbf{V}\|_F^2 \right\}. \end{aligned}$$

Rearranging terms, both problems can be written, equivalently, as

$$\operatorname{argmin}_{\mathbf{V} \in \mathcal{D}_M} \left\{ \|\mathbf{V} - \mathbf{U}\|_F^2 \right\}, \quad (10)$$

for $\mathbf{U} = \mathbf{A}_k - L^{-1}\mathbf{G}_k$ and $\mathbf{U} = -L^{-1}\mathbf{W}_k$ for the \mathbf{B}_k and \mathbf{C}_k iterations respectively. The solution of (10) is given by \mathbf{V}^* defined entrywise by (cf. Section 5.2.2 in [14])

$$\mathbf{V}^* = \frac{M}{2} \text{sign}(\mathbf{U}) \min\left(\frac{2}{M}|\mathbf{U}|, 1\right).$$

Next, we show that our notion of δ -oracle yields a δ' -approximate gradient in the sense of Equation (2.3) in [13]. Precisely, we prove that

$$\left| \text{tr}\left(\left(\tilde{D}\Phi_{[\mathbf{A}]} - D\Phi_{[\mathbf{A}]}\right)^\top (\mathbf{B} - \mathbf{C})\right) \right| \leq \delta', \quad (11)$$

for any $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathcal{D}_M$. By Hölder's inequality,

$$\left| \text{tr}\left(\left(\tilde{D}\Phi_{[\mathbf{A}]} - D\Phi_{[\mathbf{A}]}\right)^\top (\mathbf{B} - \mathbf{C})\right) \right| \leq M \left\| \tilde{D}\Phi_{[\mathbf{A}]} - D\Phi_{[\mathbf{A}]} \right\|_1,$$

and the choice $\mathbf{B} = -\mathbf{C} = \frac{M}{2} \text{sign} \left(\tilde{D}\Phi_{[\mathbf{A}]} - D\Phi_{[\mathbf{A}]} \right)$ saturates the above bound. Recall that

$$\tilde{D}\Phi_{[\mathbf{A}]} - D\Phi_{[\mathbf{A}]} = 32 \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} x^{(i)} \left(y^{(j)} \right)^\top \left(\tilde{\Pi}_{ij}^{\mathbf{A}} - \Pi_{ij}^{\mathbf{A}} \right), \quad (12)$$

where $\left\| \tilde{\Pi}^{\mathbf{A}} - \Pi^{\mathbf{A}} \right\|_\infty < \delta$ uniformly in $\mathbf{A} \in \mathcal{D}_M$ by the δ -oracle such that

$$\left\| \tilde{D}\Phi_{[\mathbf{A}]} - D\Phi_{[\mathbf{A}]} \right\|_1 \leq 32 \left\| \tilde{\Pi}^{\mathbf{A}} - \Pi^{\mathbf{A}} \right\|_\infty \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} \left\| x^{(i)} \left(y^{(j)} \right)^\top \right\|_1 < 32\delta \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} \left\| x^{(i)} \left(y^{(j)} \right)^\top \right\|_1,$$

where $|\cdot|$ is applied componentwise in the above display. Combining the displayed equations yields

$$\left| \text{tr} \left(\left(\tilde{D}\Phi_{[\mathbf{A}]} - D\Phi_{[\mathbf{A}]} \right)^\top (\mathbf{B} - \mathbf{C}) \right) \right| \leq 32M\delta \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} \left\| x^{(i)} \left(y^{(j)} \right)^\top \right\|_1 = \delta',$$

proving Eq. (11).

With these preparations Theorem 2 follows from Theorem 2.2 in [13] and the discussion following its proof, noting that $\sum_{i=0}^k \frac{i+1}{2} = \frac{(k+1)(k+2)}{4}$. \square

A.5 Proof of Corollary 2

As $\mathbf{A}_k, \mathbf{B}_k$ be iterates from Algorithm 2 with $\mathbf{B}_k \in \text{int}(\mathcal{D}_M)$ such that $\mathbf{B}_k = \mathbf{A}_k - \beta_k \tilde{D}\Phi_{[\mathbf{A}_k]}$ by definition. By the triangle inequality,

$$\|D\Phi_{[\mathbf{A}_k]}\|_F \leq \|D\Phi_{[\mathbf{A}_k]} - \tilde{D}\Phi_{[\mathbf{A}_k]}\|_F + \|\tilde{D}\Phi_{[\mathbf{A}_k]}\|_F = \|D\Phi_{[\mathbf{A}_k]} - \tilde{D}\Phi_{[\mathbf{A}_k]}\|_F + \|\beta_k^{-1} (\mathbf{B}_k - \mathbf{A}_k)\|_F.$$

It remains to bound $\|D\Phi_{[\mathbf{A}_k]} - \tilde{D}\Phi_{[\mathbf{A}_k]}\|_F$ using the δ -oracle. By (12),

$$\begin{aligned} \|D\Phi_{[\mathbf{A}]} - \tilde{D}\Phi_{[\mathbf{A}]}\|_F &= 32 \left\| \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} x^{(i)} \left(y^{(j)} \right)^\top \left(\tilde{\Pi}_{ij}^{\mathbf{A}} - \Pi_{ij}^{\mathbf{A}} \right) \right\|_F \\ &\leq 32 \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} \left| \tilde{\Pi}_{ij}^{\mathbf{A}} - \Pi_{ij}^{\mathbf{A}} \right| \left\| x^{(i)} \left(y^{(j)} \right)^\top \right\|_F \\ &\leq 32 \|\tilde{\Pi}^{\mathbf{A}} - \Pi^{\mathbf{A}}\|_\infty \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} \|x^{(i)}\| \|y^{(j)}\| \\ &< 32\delta \sum_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} \|x^{(i)}\| \|y^{(j)}\| \end{aligned}$$

\square

B Sharpness of variance bound from Lemma 4

Let $\mu_0 = \frac{1}{2}(\delta_0 + \delta_a)$ and $\mu_1 = \frac{1}{2}(\delta_0 + \delta_b)$ for $a \in \mathbb{R}^{d_0}$ and $b \in \mathbb{R}^{d_1}$. In this case, any coupling $\pi \in \Pi(\mu_0, \mu_1)$ is of the form $\pi_{00}\delta_{(0,0)} + \pi_{0b}\delta_{(0,b)} + \pi_{a0}\delta_{(a,0)} + \pi_{ab}\delta_{(a,b)}$ with the constraint that $\pi_{00} = \pi_{ab}$ and $\pi_{0b} = \pi_{a0} = \frac{1}{2} - \pi_{ab}$. For any $\mathbf{A} \in \mathcal{D}_M$, $\text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1)$ is given by

$$\begin{aligned} &\inf_{\pi \in \Pi(\mu_0, \mu_1)} \left\{ \int -4\|x\|^2 \|y\|^2 - 32x^\top \mathbf{A} y \, d\pi(x, y) + \varepsilon \text{D}_{\text{KL}}(\pi \| \mu_0 \otimes \mu_1) \right\} \\ &= \inf_{\pi_{ab} \in (0, 1/2)} \left\{ -\pi_{ab}(4\|a\|^2 \|b\|^2 + 32a^\top \mathbf{A} b) + 2\varepsilon \pi_{ab} \log(4\pi_{ab}) + (1 - 2\pi_{ab}) \varepsilon \log(2 - 4\pi_{ab}) \right\}, \end{aligned}$$

the objective is a sum of convex functions and the first-order optimality condition reads

$$4\|a\|^2\|b\|^2 + 32a^\top \mathbf{A}b = 2\varepsilon \log(4\pi_{ab}) - 2\varepsilon \log(2 - 4\pi_{ab}) \iff \pi_{ab} = \frac{e^z}{2(1 + e^z)},$$

for $z = (2\|a\|^2\|b\|^2 + 16a^\top \mathbf{A}b) / \varepsilon$. Let π^* be the corresponding EOT coupling for $\text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1)$. For any $\mathbf{C} \in \mathbb{R}^{d_0 \times d_1}$,

$$\text{Var}_{\pi^*}[X^\top \mathbf{C}Y] = \pi_{ab}^*(1 - \pi_{ab}^*)(a^\top \mathbf{C}b)^2 \leq \pi_{ab}^*(1 - \pi_{ab}^*)\|\mathbf{C}\|_F^2\|a\|^2\|b\|^2,$$

with equality for $\mathbf{C} = Cab^\top$ with $C \in \mathbb{R}$. Hence,

$$\sup_{\|\mathbf{C}\|_F=1} \{\text{Var}_{\pi^*}[X^\top \mathbf{C}Y]\} = \pi_{ab}^*(1 - \pi_{ab}^*)\|a\|^2\|b\|^2,$$

which can be made arbitrarily close to $\frac{1}{4}\|a\|^2\|b\|^2$ for fixed a, b by choosing $\mathbf{A} \in \mathcal{D}_M$ and $\varepsilon > 0$ as to make z sufficiently large. On the other hand, $\sqrt{M_4(\mu_0)M_4(\mu_1)} = \frac{1}{2}\|a\|^2\|b\|^2$, such that the variance bound in Lemma 4 is tight up to a constant factor.

C Sinkhorn's Algorithm as an inexact oracle

Given $\mu_0 = \sum_{i=1}^{N_0} a_i \delta_{x^{(i)}} \in \mathcal{P}(\mathbb{R}^{d_0})$ and $\mu_1 = \sum_{j=1}^{N_1} b_j \delta_{x^{(j)}} \in \mathcal{P}(\mathbb{R}^{d_1})$, let a, b denote the corresponding (positive) probability vectors. Fix an underlying cost function $c : \mathbb{R}^{d_0} \times \mathbb{R}^{d_1} \rightarrow \mathbb{R}$ and $\varepsilon > 0$, and let $\mathbf{K} \in \mathbb{R}^{N_0 \times N_1}$ with $\mathbf{K}_{ij} = e^{-\frac{c(x^{(i)}, y^{(j)})}{\varepsilon}}$. Consider the standard implementation of Sinkhorn's algorithm (cf. e.g. [12, 19]).

Algorithm 3 Sinkhorn Algorithm

- 1: Fix a threshold γ and a maximum iteration number k_{\max} .
 - 2: $v_0 \leftarrow \mathbb{1}_{N_1}/N_1$
 - 3: $k \leftarrow 1$
 - 4: **repeat**
 - 5: $u_k \leftarrow a/(\mathbf{K}v_{k-1})$
 - 6: $v_k \leftarrow b/(\mathbf{K}^\top u_k)$
 - 7: $\mathbf{\Pi}^k \leftarrow \text{diag}(u_k)\mathbf{K}\text{diag}(v_k)$
 - 8: $k \leftarrow k + 1$
 - 9: **until** $\|\mathbf{\Pi}^k \mathbb{1}_{N_1} - a\|_2 < \gamma$ or $k > k_{\max}$
 - 10: **return** $\mathbf{\Pi}^k$
-

In Algorithm 3, the division of two vectors is understood componentwise. The stopping condition is based only on one of the marginal constraints as $\mathbb{1}_{N_0}^\top \mathbf{\Pi}^k = b^\top$ by construction.

The following definitions enable describing the convergence properties of Algorithm 3; we follow the approach of [20] with only minor modifications. Let \mathbb{R}_+^d denote the set of vectors with positive entries and, for $x, y \in \mathbb{R}_+^d$ let

$$d_H(x, y) = \log \max_{1 \leq i, j \leq d} \frac{x_i y_j}{y_i x_j},$$

denote Hilbert's projective metric³ on \mathbb{R}_+^d . By definition,

$$d_H(x, y) = d_H(x/y, \mathbb{1}_d), \tag{13}$$

for any $x, y \in \mathbb{R}_+^d$ and, setting $x = e^w, y = e^z$ componentwise,

$$\begin{aligned} d_H(x, y) &= \log \max_{1 \leq i, j \leq d} e^{w_i + z_j - w_j - z_i}, \\ &= \max_{1 \leq i, j \leq d} w_i + z_j - w_j - z_i, \\ &= \max_{1 \leq i \leq d} (\log x_i - \log y_i) - \min_{1 \leq i \leq d} (\log x_i - \log y_i), \\ &= \max_{1 \leq i \leq d} \log \left(\frac{x_i}{y_i} \right) - \min_{1 \leq i \leq d} \log \left(\frac{x_i}{y_i} \right). \end{aligned} \tag{14}$$

³ $d_H(x, y) = 0$ if and only if $x = \alpha y$ for $\alpha > 0$, d_H is symmetric and satisfies the triangle inequality.

It was proved in [4, 31] that multiplication with a positive matrix is a strict contraction w.r.t. d_H . Precisely,

$$d_H(\mathbf{A}x, \mathbf{A}y) \leq \lambda(\mathbf{A})d_H(x, y), \quad (15)$$

for any $\mathbf{A} \in \mathbb{R}_+^{d' \times d}$ and $x, y \in \mathbb{R}_+^d$, where

$$\lambda(\mathbf{A}) = \frac{\sqrt{\eta(\mathbf{A})} - 1}{\sqrt{\eta(\mathbf{A})} + 1} < 1, \quad \eta(\mathbf{A}) = \max_{\substack{1 \leq i, j \leq d' \\ 1 \leq k, l \leq d}} \frac{\mathbf{A}_{ik}\mathbf{A}_{jl}}{\mathbf{A}_{jk}\mathbf{A}_{il}}.$$

Let

$$E = \{\mathbf{A} \in \mathbb{R}_+^{N_0 \times N_1} : \mathbf{A} = \text{diag}(x)\mathbf{K}\text{diag}(y) \text{ for some } x \in \mathbb{R}_+^{N_0}, y \in \mathbb{R}_+^{N_1}\},$$

and observe that if $\mathbf{A}, \mathbf{B} \in E$, there exists $x_{\mathbf{A}, \mathbf{B}} \in \mathbb{R}_+^{N_0}, y_{\mathbf{A}, \mathbf{B}} \in \mathbb{R}_+^{N_1}$ for which $\mathbf{A} = \text{diag}(x_{\mathbf{A}, \mathbf{B}})\mathbf{B}\text{diag}(y_{\mathbf{A}, \mathbf{B}})$. In this setting, let $d : E \times E \mapsto [0, \infty)$ be such that

$$d(\mathbf{A}, \mathbf{B}) = d_H(x_{\mathbf{A}, \mathbf{B}}, \mathbf{1}_{N_0}) + d_H(y_{\mathbf{A}, \mathbf{B}}, \mathbf{1}_{N_1}),$$

then d is a metric on E . As the EOT coupling $\mathbf{\Pi}^*$ satisfies

$$\frac{\mathbf{\Pi}_{ij}^*}{a_i b_j} = e^{\frac{\varphi(x^{(i)}) + \psi(y^{(j)}) - c(x^{(i)}, y^{(j)})}{\varepsilon}},$$

where (φ, ψ) is any pair of EOT potentials, $\mathbf{\Pi}^* = \text{diag}(u^*)\mathbf{K}\text{diag}(v^*) \in E$ for

$$u_i^* = a_i e^{\frac{\varphi(x^{(i)})}{\varepsilon}}, \quad v_j^* = b_j e^{\frac{\psi(y^{(j)})}{\varepsilon}}.$$

Note that $u^* = a/(\mathbf{K}v^*)$ and $v^* = b/(\mathbf{K}^\top u^*)$.

In the sequel, we analyze the convergence of $\mathbf{\Pi}^k$ to $\mathbf{\Pi}^*$ under d . The following result translates bounds on $d(\mathbf{\Pi}^k, \mathbf{\Pi}^*)$ to bounds on $\|\mathbf{\Pi}^k - \mathbf{\Pi}^*\|_\infty$.

Lemma 5. Fix $\delta > 0$. If $d(\mathbf{\Pi}^k, \mathbf{\Pi}^*) \leq \delta$, it follows that $\|\mathbf{\Pi}^k - \mathbf{\Pi}^*\|_\infty \leq e^\delta - 1$.

Proof. By Lemma 3 in [20], whenever $d(\mathbf{\Pi}^k, \mathbf{\Pi}^*) \leq \delta$ it holds that

$$e^{-\delta} - 1 \leq \frac{\mathbf{\Pi}_{ij}^*}{\mathbf{\Pi}_{ij}^k} - 1 \leq e^\delta - 1,$$

for every $1 \leq i \leq N_0, 1 \leq j \leq N_1$. As such,

$$|\mathbf{\Pi}_{ij}^* - \mathbf{\Pi}_{ij}^k| \leq \mathbf{\Pi}_{ij}^k ((1 - e^{-\delta}) \vee (e^\delta - 1)) \leq (1 - e^{-\delta}) \vee (e^\delta - 1) = e^\delta - 1,$$

yielding $\|\mathbf{\Pi}^* - \mathbf{\Pi}^k\|_\infty \leq e^\delta - 1$. \square

Towards bounding $d(\mathbf{\Pi}^k, \mathbf{\Pi}^*)$, we first show that the iterates u_k, v_k defined in Algorithm 3 converge to u^*, v^* under d_H .

Lemma 6. Let u_k, v_k be iterates generated by Algorithm 3. Then,

$$\begin{aligned} d_H(u_k, u^*) &\leq \lambda(\mathbf{K})^{2(k-1)} d_H(u_1, u^*) \\ d_H(v_k, v^*) &\leq \lambda(\mathbf{K})^{2k} d_H(v_0, v^*). \end{aligned}$$

In particular, $d_H(u_k, u^*) \rightarrow 0, d_H(v_k, v^*) \rightarrow 0$ as $k \rightarrow \infty$.

Proof. The second claim follows from the first and the fact that $\lambda(\mathbf{K}) < 1$. To prove the first claim, we have, by definition,

$$\begin{aligned} d_H(u_{k+1}, u^*) &= d_H\left(\frac{a}{\mathbf{K}v_k}, \frac{a}{\mathbf{K}v^*}\right) = d_H(\mathbf{K}v_k, \mathbf{K}v^*) \leq \lambda(\mathbf{K})d_H(v_k, v^*), \\ d_H(v_k, v^*) &= d_H\left(\frac{b}{\mathbf{K}^\top u_k}, \frac{b}{\mathbf{K}^\top u^*}\right) = d_H(\mathbf{K}^\top u_k, \mathbf{K}^\top u^*) \leq \lambda(\mathbf{K})d_H(u_k, u^*), \end{aligned} \quad (16)$$

as $\lambda(\mathbf{K}) = \lambda(\mathbf{K}^\top)$. Thus, $d_H(u_{k+1}, u^*) \leq \lambda(\mathbf{K})^2 d_H(u_k, u^*)$ and $d_H(v_{k+1}, v^*) \leq \lambda(\mathbf{K})^2 d_H(v_k, v^*)$. The conclusion follows by applying these bounds repeatedly. \square

Next, we bound the progress of the iterates $\mathbf{\Pi}^k$ to $\mathbf{\Pi}^*$ under d in terms of $d_H(u_k, u^*)$, $d_H(v_k, v^*)$.

Lemma 7. *In the setting of Lemma 6, we have that*

$$d(\mathbf{\Pi}^k, \mathbf{\Pi}^*) = d_H(u_k, u^*) + d_H(v_k, v^*) \leq \lambda(\mathbf{K})^{2(k-1)} d_H(u_1, u^*) + \lambda(\mathbf{K})^{2k} d_H(v_0, v^*).$$

Proof. The second inequality follows from the first by Lemma 6. By construction,

$$\mathbf{\Pi}^{k+j} = \text{diag} \left(\frac{u_j}{u_k} \right) \mathbf{\Pi}^k \text{diag} \left(\frac{v_j}{v_k} \right),$$

such that

$$\begin{aligned} d(\mathbf{\Pi}^k, \mathbf{\Pi}^{k+j}) &= d_H \left(\frac{u_j}{u_k}, \mathbf{1}_{N_0} \right) + d_H \left(\frac{v_j}{v_k}, \mathbf{1}_{N_1} \right), \\ &= d_H(u_j, u_k) + d_H(v_j, v_k), \end{aligned}$$

taking the limit as $j \rightarrow \infty$ and applying Lemma 6 yields

$$d(\mathbf{\Pi}^k, \mathbf{\Pi}^*) = d_H(u_k, u^*) + d_H(v_k, v^*),$$

where we have also used the fact that $\lim_{j \rightarrow \infty} d(\mathbf{\Pi}^k, \mathbf{\Pi}^{k+j}) = d(\mathbf{\Pi}^k, \mathbf{\Pi}^*)$ as follows from [20, p. 731] and [35]. \square

As u^* and v^* are *a priori* unknown, we now bound $d(u^k, u^*)$ in terms of $d(a, u_k \odot \mathbf{K}v_k)$, which is a measure of how much $\mathbf{\Pi}^k$ violates the marginal constraint,⁴ and another analogous term.

Lemma 8. *In the setting of Lemma 6, we have that*

$$d_H(u_k, u^*) \leq \frac{d_H(a, u_k \odot \mathbf{K}v_k)}{1 - \lambda(\mathbf{K})^2}, \quad d_H(v_k, v^*) \leq \frac{d_H(a, v_k \odot \mathbf{K}^\top u_{k+1})}{1 - \lambda(\mathbf{K})^2}.$$

Proof. By construction, we have that

$$\begin{aligned} d_H(u_k, u^*) &\leq d_H(u_{k+1}, u_k) + d_H(u_{k+1}, u^*), \\ &\leq d_H \left(\frac{a}{\mathbf{K}v_k}, u_k \right) + \lambda(\mathbf{K})^2 d_H(u_k, u^*), \\ &= d_H(a, u_k \odot \mathbf{K}v_k) + \lambda(\mathbf{K})^2 d_H(u_k, u^*), \end{aligned}$$

where we have applied the triangle inequality and (16). The claimed result for v_k follows from the same argument. \square

Combined, Lemmas 7 and 8 provide an explicit bound on the total number of iterations required to ensure that $d(\mathbf{\Pi}^k, \mathbf{\Pi}^*)$ achieves a given precision.

Proposition 3. *Let $\mathbf{\Pi}^k$ be given by Algorithm 3 and fix $\delta > 0$. Then, for every*

$$k \geq 1 + \frac{1}{2 \log(\lambda(\mathbf{K}))} \log \left(\frac{\delta(1 - \lambda(\mathbf{K})^2)}{d_H(a, u_1 \odot \mathbf{K}v_1) + \lambda(\mathbf{K})^2 d_H(b, v_0 \odot \mathbf{K}^\top u_1)} \right),$$

$$d(\mathbf{\Pi}^k, \mathbf{\Pi}^*) \leq \delta.$$

Proof. It follows from Lemma 7 and Lemma 8 that

$$d(\mathbf{\Pi}^k, \mathbf{\Pi}^*) \leq \frac{\lambda(\mathbf{K})^{2(k-1)} d_H(a, u_1 \odot \mathbf{K}v_1) + \lambda(\mathbf{K})^{2k} d_H(b, v_0 \odot \mathbf{K}^\top u_1)}{1 - \lambda(\mathbf{K})^2}.$$

The upper bound on the number of iterations required to achieve $d(\mathbf{\Pi}^k, \mathbf{\Pi}^*) \leq \delta$ then follows from basic algebra. \square

⁴ $\mathbf{\Pi}^k \mathbf{1}_{N_1} = u_k \odot \mathbf{K}v_k$, where \odot denotes elementwise multiplication.

Now, we demonstrate why the termination condition based on the 2-norm (or an equivalent condition based on the 1-norm) endows us with a δ -oracle approximation and provide a bound on the number of iterations required to achieve it. Theorem 1 in [18] proves that there exists $\bar{k} \leq 1 + \frac{R}{\delta}$ satisfying

$$\|u_{\bar{k}} \odot \mathbf{K}v_{\bar{k}} - a\|_1 + \|v_{\bar{k}} \odot \mathbf{K}^\top u_{\bar{k}+1} - b\|_1 \leq \delta,$$

for $R = -2 \log \left(e^{-\|\mathbf{C}\|_\infty/\varepsilon} \min_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} a_i \wedge b_j \right)$. This gives a bound on the maximal number of iterations to achieve the 2-norm termination condition via the standard inequality $\|\cdot\|_2 \leq \|\cdot\|_1$.

We now bound d_H in terms of the Euclidean distance as to control $d(\Pi^k, \Pi^*)$ by $\|u_k \odot \mathbf{K}v_k - a\|_2$.

Lemma 9. *Let $r, s \in \mathbb{R}_+^d$ be arbitrary, then*

$$d_H(s, r) \leq (r_{i_*}^{-1} + s_{i_*}^{-1}) \|r - s\|_2,$$

where $i_* \in \operatorname{argmax}_{1 \leq i \leq d} r_i^{-1}(s_i - r_i)$ and $i_* \in \operatorname{argmin}_{1 \leq i \leq d} s_i^{-1}(s_i - r_i)$.

Proof. We have by (14) that

$$d_H(s, r) = \max_{1 \leq i \leq N_0} \log \left(\frac{s_i}{r_i} \right) - \min_{1 \leq i \leq N_0} \log \left(\frac{s_i}{r_i} \right).$$

Observe that

$$1 - \frac{r_i}{s_i} \leq \log \left(\frac{s_i}{r_i} \right) \leq \frac{s_i}{r_i} - 1,$$

using the inequalities $\frac{x}{1+x} \leq \log(1+x) \leq x$ for $x > -1$. Whence,

$$\begin{aligned} d_H(s, r) &\leq \max_{1 \leq i \leq N_0} r_i^{-1}(s_i - r_i) - \min_{1 \leq i \leq N_0} s_i^{-1}(s_i - r_i) \\ &= r_{i_*}^{-1}(s_{i_*} - r_{i_*}) - s_{i_*}^{-1}(s_{i_*} - r_{i_*}) \\ &\leq (r_{i_*}^{-1} + s_{i_*}^{-1}) \|s - r\|_2. \end{aligned}$$

□

By combining Lemmas 7 and 9 we arrive at the desired result.

Proposition 4. *Let $\underline{a} = \min_{1 \leq i \leq N_0} a_i$. In the setting of Lemma 9, the iterates Π^k generated by Algorithm 3 with the threshold $\underline{a} \geq \gamma > 0$ satisfy*

$$d(\Pi^k, \Pi^*) \leq \frac{E_k \|a - u_k \odot \mathbf{K}v_k\|_2}{1 - \lambda(\mathbf{K})},$$

where E_k denotes the constant from Lemma 9. Hence, the 2-norm termination criterion in Algorithm 3 is satisfied in \bar{k} iterations for some $\bar{k} \leq 1 + \frac{R}{\gamma}$ and

$$d(\Pi^{\bar{k}}, \Pi^*) \leq \frac{E_{\bar{k}} \gamma}{1 - \lambda(\mathbf{K})} \leq \frac{\gamma}{1 - \lambda(\mathbf{K})} (\underline{a}^{-1} + (\underline{a} - \gamma)^{-1}).$$

Proof. The first claim follows directly from Lemmas 7 and 9 together with the fact that $d_H(v_k, v^*) \leq \lambda(\mathbf{K}) d_H(u_k, u^*)$ (see (16)).

It is clear from the discussion preceding Lemma 9 that $\|u_{\bar{k}} \odot \mathbf{K}v_{\bar{k}} - a\|_2 \leq \gamma$ for some $\bar{k} \leq 1 + \frac{R}{\gamma}$, which corresponds to the 2-norm termination condition for Algorithm 3. To see that $E_{\bar{k}} \leq \underline{a}^{-1} + (\underline{a} - \gamma)^{-1}$, let $w^k = u_k \odot \mathbf{K}v_k$ and observe that $\|a - w^k\|_\infty \leq \|a - w^k\|_2 \leq \gamma$. Hence, for any index $1 \leq i \leq N_0$, $a_i - \gamma \leq w_i^k$ and, as such, $(w_i^k)^{-1} \leq \frac{1}{a_i - \gamma} \leq \frac{1}{\underline{a} - \gamma}$. □

The proof of Proposition 2 then follows by combining Propositions 3 and Proposition 4. The maximal number of iterations for Algorithm 3 to output a δ -oracle approximation of the EOT coupling

Π^A is thus

$$\tilde{k} = \min \left\{ 1 + \frac{1}{2 \log(\lambda(\mathbf{K}))} \log \left(\frac{\delta(1 - \lambda(\mathbf{K})^2)}{\mathbf{d}_H(a, u_1 \odot \mathbf{K} v_1) + \lambda(\mathbf{K})^2 \mathbf{d}_H(b, v_0 \odot \mathbf{K}^\top u_1)} \right), \right. \\ \left. 1 - 2\delta^{-1} \log \left(e^{-\|\mathbf{C}\|_\infty/\varepsilon} \min_{\substack{1 \leq i \leq N_0 \\ 1 \leq j \leq N_1}} a_i \wedge b_j \right) \right\}. \quad (17)$$

D Convergence of Algorithm 2

In what follows, we slightly adapt the proof of Theorem 2 in [21] to conform to the inexact setting.

We first clarify that they treat the composite problem

$$\inf_{x \in \mathbb{R}^d} f(x) + g(x) + \mathcal{Q}(x),$$

where f is L' -smooth and non-convex, g is L'' -smooth and convex, and \mathcal{Q} is non-smooth and convex with a bounded domain. Hence $f + g$ is $L = L' + L''$ smooth and possibly non-convex.

Our problem conforms to this setting (up to vectorization) with $f = \text{OT}_{(\cdot), \varepsilon}(\mu_0, \mu_1)$, $g = 32\|\cdot\|_F^2$, and $\mathcal{Q} = \mathcal{I}_{\mathcal{D}_M}$, the indicator function of the set \mathcal{D}_M , defined by

$$\mathcal{I}_{\mathcal{D}_M}(\mathbf{A}) = \begin{cases} 0, & \text{if } \mathbf{A} \in \mathcal{D}_M, \\ +\infty, & \text{otherwise.} \end{cases}$$

When Φ is convex, we set $f = 0$ and $g = \Phi$ hence $L' = 0$, $L = L''$.

As Φ is L -smooth, by Lemma 5 in [21],

$$\Phi(\mathbf{B}_k) \leq \Phi(\mathbf{A}_k) + \text{tr} \left(D\Phi_{[\mathbf{A}_k]}^\top (\mathbf{B}_k - \mathbf{A}_k) \right) + \frac{L}{2} \|\mathbf{B}_k - \mathbf{A}_k\|_F^2, \quad (18)$$

and for any $\mathbf{H} \in \mathbb{R}^{d_0 \times d_1}$, letting L' denote the Lipschitz constant of $\text{OT}_{(\cdot), \varepsilon}(\mu_0, \mu_1)$, the same result gives

$$\begin{aligned} & \Phi(\mathbf{A}_k) - ((1 - \tau_k)\Phi(\mathbf{B}_{k-1}) + \tau_k\Phi(\mathbf{H})) \\ &= \tau_k (\Phi(\mathbf{A}_k) - \Phi(\mathbf{H})) + (1 - \tau_k) (\Phi(\mathbf{A}_k) - \Phi(\mathbf{B}_{k-1})) \\ &\leq \tau_k \left(\text{tr} \left(D\Phi_{[\mathbf{A}_k]}^\top (\mathbf{A}_k - \mathbf{H}) \right) + \frac{L'}{2} \|\mathbf{H} - \mathbf{A}_k\|_F^2 \right) \\ &+ (1 - \tau_k) \left(\text{tr} \left(D\Phi_{[\mathbf{A}_k]}^\top (\mathbf{A}_k - \mathbf{B}_{k-1}) \right) + \frac{L'}{2} \|\mathbf{B}_{k-1} - \mathbf{A}_k\|_F^2 \right) \\ &= \text{tr} \left(D\Phi_{[\mathbf{A}_k]}^\top (\mathbf{A}_k - \tau_k\mathbf{H} - (1 - \tau_k)\mathbf{B}_{k-1}) \right) \\ &+ \frac{L'\tau_k}{2} \|\mathbf{H} - \mathbf{A}_k\|_F^2 + \frac{L'(1 - \tau_k)}{2} \underbrace{\|\mathbf{B}_{k-1} - \mathbf{A}_k\|_F^2}_{\tau_k^2 \|\mathbf{B}_{k-1} - \mathbf{C}_{k-1}\|_F^2}, \end{aligned} \quad (19)$$

recalling the update $\mathbf{A}_k = \tau_k \mathbf{C}_{k-1} + (1 - \tau_k) \mathbf{B}_{k-1}$.

Denote the subdifferential of $\mathcal{I}_{\mathcal{D}_M}$ at $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$ by

$$\partial \mathcal{I}_{\mathcal{D}_M}(\mathbf{A}) := \left\{ \mathbf{P} \in \mathbb{R}^{d_0 \times d_1} : \mathcal{I}_{\mathcal{D}_M}(\mathbf{X}) - \mathcal{I}_{\mathcal{D}_M}(\mathbf{A}) \geq \text{tr}(\mathbf{P}^\top (\mathbf{X} - \mathbf{A})), \text{ for every } \mathbf{X} \in \mathbb{R}^{d_0 \times d_1} \right\}.$$

As \mathbf{C}_k is optimal for the problem $\text{argmin}_{\mathbf{V} \in \mathbb{R}^{d_0 \times d_1}} \left\{ \frac{1}{2\gamma_k} \|\mathbf{V} - (\mathbf{C}_{k-1} - \gamma_k \mathbf{G}_k)\|_F^2 + \mathcal{I}_{\mathcal{D}_M}(\mathbf{V}) \right\}$, there exists $\mathbf{P} \in \partial \mathcal{I}_{\mathcal{D}_M}(\mathbf{C}_k)$ for which $\mathbf{G}_k + \mathbf{P} + \frac{1}{\gamma_k} (\mathbf{C}_k - \mathbf{C}_{k+1}) = 0$ (see Theorem 23.8, Theorem 25.1, and p. 264 in [30]). Thus, for any $\mathbf{U} \in \mathbb{R}^{d_0 \times d_1}$,

$$\begin{aligned} \text{tr}((\mathbf{G}_k + \mathbf{P})^\top (\mathbf{C}_k - \mathbf{U})) &= \frac{1}{\gamma_k} \text{tr}((\mathbf{C}_k - \mathbf{C}_{k-1})^\top (\mathbf{U} - \mathbf{C}_k)) \\ &= \frac{1}{2\gamma_k} \left(\|\mathbf{C}_{k-1} - \mathbf{U}\|_F^2 - \|\mathbf{C}_k - \mathbf{U}\|_F^2 - \|\mathbf{C}_k - \mathbf{C}_{k-1}\|_F^2 \right), \end{aligned}$$

where the final line follows from some simple algebra. As $\mathbf{P} \in \partial \mathcal{I}_{\mathcal{D}_M}(\mathbf{C}_k)$, $\text{tr}(\mathbf{P}^\top(\mathbf{C}_k - \mathbf{U})) \geq \mathcal{I}_{\mathcal{D}_M}(\mathbf{C}_k) - \mathcal{I}_{\mathcal{D}_M}(\mathbf{U}) = -\mathcal{I}_{\mathcal{D}_M}(\mathbf{U})$, whence

$$\text{tr}(\mathbf{G}_k^\top(\mathbf{C}_k - \mathbf{U})) \leq \mathcal{I}_{\mathcal{D}_M}(\mathbf{U}) + \frac{1}{2\gamma_k} (\|\mathbf{C}_{k-1} - \mathbf{U}\|_F^2 - \|\mathbf{C}_k - \mathbf{U}\|_F^2 - \|\mathbf{C}_k - \mathbf{C}_{k-1}\|_F^2). \quad (20)$$

By the same steps applied to the other subproblem with \mathbf{B}_k and \mathbf{A}_k taking the place of \mathbf{C}_k and \mathbf{C}_{k-1} respectively,

$$\text{tr}(\mathbf{G}_k^\top(\mathbf{B}_k - \mathbf{U})) \leq \mathcal{I}_{\mathcal{D}_M}(\mathbf{U}) + \frac{1}{2\beta_k} (\|\mathbf{A}_k - \mathbf{U}\|_F^2 - \|\mathbf{B}_k - \mathbf{U}\|_F^2 - \|\mathbf{B}_k - \mathbf{A}_k\|_F^2).$$

Setting $\mathbf{U} = \tau_k \mathbf{C}_k + (1 - \tau_k) \mathbf{B}_{k-1} \in \mathcal{D}_M$ (by convexity) in the previous display, bounding $-\|\mathbf{B}_k - \mathbf{U}\|_F^2$ above by 0, and recalling that $\mathbf{A}_k = \tau_k \mathbf{C}_{k-1} + (1 - \tau_k) \mathbf{B}_{k-1}$ such that $\mathbf{A}_k - \mathbf{U} = \tau_k(\mathbf{C}_{k-1} - \mathbf{C}_k)$,

$$\text{tr}(\mathbf{G}_k^\top(\mathbf{B}_k - \tau_k \mathbf{C}_k + (1 - \tau_k) \mathbf{B}_{k-1})) \leq \frac{1}{2\beta_k} (\tau_k^2 \|\mathbf{C}_k - \mathbf{C}_{k-1}\|_F^2 - \|\mathbf{B}_k - \mathbf{A}_k\|_F^2).$$

Combining with (20) upon scaling by τ_k ,

$$\begin{aligned} \text{tr}(\mathbf{G}_k^\top(\mathbf{B}_k - \tau_k \mathbf{U} + (1 - \tau_k) \mathbf{B}_{k-1})) &\leq \tau_k \mathcal{I}_{\mathcal{D}_M}(\mathbf{U}) + \frac{1}{2\beta_k} (\tau_k^2 \|\mathbf{C}_k - \mathbf{C}_{k-1}\|_F^2 - \|\mathbf{B}_k - \mathbf{A}_k\|_F^2), \\ &\quad + \frac{\tau_k}{2\gamma_k} (\|\mathbf{C}_{k-1} - \mathbf{U}\|_F^2 - \|\mathbf{C}_k - \mathbf{U}\|_F^2 - \|\mathbf{C}_k - \mathbf{C}_{k-1}\|_F^2), \end{aligned}$$

by the choice of $\tau_k, \beta_k, \gamma_k$, we have that $\frac{\tau_k^2}{\beta_k} - \frac{\tau_k}{\gamma_k} \leq 0$ such that

$$\begin{aligned} \text{tr}(\mathbf{G}_k^\top(\mathbf{B}_k - \tau_k \mathbf{U} + (1 - \tau_k) \mathbf{B}_{k-1})) &\leq \tau_k \mathcal{I}_{\mathcal{D}_M}(\mathbf{U}) + \frac{\tau_k}{2\gamma_k} (\|\mathbf{C}_{k-1} - \mathbf{U}\|_F^2 - \|\mathbf{C}_k - \mathbf{U}\|_F^2) \\ &\quad - \frac{1}{2\beta_k} \|\mathbf{B}_k - \mathbf{A}_k\|_F^2. \end{aligned}$$

Combining the equation above with Eq. (18) and Eq. (19) and setting $\mathbf{H} = \mathbf{U} \in \mathcal{D}_M$ (otherwise the bound is vacuous),

$$\begin{aligned} \Phi(\mathbf{B}_k) - \Phi(\mathbf{H}) &\leq (1 - \tau_k) (\Phi(\mathbf{B}_{k-1}) - \Phi(\mathbf{H})) + \text{tr} \left(D\Phi_{[\mathbf{A}_k]}^\top(\mathbf{B}_k - \tau_k \mathbf{H} - (1 - \tau_k) \mathbf{B}_{k-1}) \right) \\ &\quad + \frac{L' \tau_k}{2} \|\mathbf{H} - \mathbf{A}_k\|_F^2 + \frac{L'(1 - \tau_k)}{2} \tau_k^2 \|\mathbf{B}_{k-1} - \mathbf{C}_{k-1}\|_F^2 + \frac{L}{2} \|\mathbf{B}_k - \mathbf{A}_k\|_F^2 \\ &\leq (1 - \tau_k) (\Phi(\mathbf{B}_{k-1}) - \Phi(\mathbf{H})) + \delta' + \frac{\tau_k}{2\gamma_k} (\|\mathbf{C}_{k-1} - \mathbf{H}\|_F^2 - \|\mathbf{C}_k - \mathbf{H}\|_F^2) \\ &\quad + \frac{L' \tau_k}{2} \|\mathbf{H} - \mathbf{A}_k\|_F^2 + \frac{L'(1 - \tau_k)}{2} \tau_k^2 \|\mathbf{B}_{k-1} - \mathbf{C}_{k-1}\|_F^2 + \left(\frac{L}{2} - \frac{1}{2\beta_k} \right) \|\mathbf{B}_k - \mathbf{A}_k\|_F^2, \end{aligned}$$

where the inequality follows from the δ -oracle which implies the bound (cf. (11))

$$\sup_{\mathbf{Y}, \mathbf{Z} \in \mathcal{D}_M} \{ |\text{tr}(\mathbf{G}_k - D\Phi_{[\mathbf{A}_k]}^\top(\mathbf{Y} - \mathbf{Z}))| \} \leq \delta',$$

observing that $\mathbf{B}_k, \tau_k \mathbf{H} + (1 - \tau_k) \mathbf{B}_{k-1} \in \mathcal{D}_M$ by convexity ($\tau_k \in (0, 1]$).

Applying Lemma 1 in [21] yields, for $A_i = \frac{2}{i(i+1)}$,

$$\begin{aligned} \frac{\Phi(\mathbf{B}_k) - \Phi(\mathbf{H})}{A_k} &\leq \sum_{i=1}^k A_i^{-1} \left(\delta' + \frac{\tau_i}{2\gamma_i} (\|\mathbf{C}_{i-1} - \mathbf{H}\|_F^2 - \|\mathbf{C}_i - \mathbf{H}\|_F^2) + \frac{L' \tau_i}{2} \|\mathbf{H} - \mathbf{A}_i\|^2 \right. \\ &\quad \left. + \frac{L'(1 - \tau_i)}{2} \tau_i^2 \|\mathbf{B}_{i-1} - \mathbf{C}_{i-1}\|_F^2 + \left(\frac{L}{2} - \frac{1}{2\beta_i} \right) \|\mathbf{B}_i - \mathbf{A}_i\|^2 \right) \\ &\leq \frac{\|\mathbf{C}_0 - \mathbf{H}\|_F^2}{2\gamma_1} + \sum_{i=1}^k A_i^{-1} \left(\delta' + \frac{L' \tau_i}{2} \|\mathbf{H} - \mathbf{A}_i\|^2 \right. \\ &\quad \left. + \frac{L'(1 - \tau_i)}{2} \tau_i^2 \|\mathbf{B}_{i-1} - \mathbf{C}_{i-1}\|_F^2 + \left(\frac{L}{2} - \frac{1}{2\beta_i} \right) \|\mathbf{B}_i - \mathbf{A}_i\|^2 \right). \end{aligned}$$

By convexity of $\|\cdot\|_F^2$,

$$\begin{aligned}
& \|\mathbf{H} - \mathbf{A}_i\|_F^2 + \tau_i(1 - \tau_i)\|\mathbf{B}_{i-1} - \mathbf{C}_{i-1}\|_F^2 \\
& \leq 2\left(\|\mathbf{H}\|_F^2 + \|\mathbf{A}_i\|_F^2 + \tau_i(1 - \tau_i)\left(\|\mathbf{B}_{i-1}\|_F^2 + \|\mathbf{C}_{i-1}\|_F^2\right)\right) \\
& \leq 2\left(\|\mathbf{H}\|_F^2 + (1 - \tau_i)\|\mathbf{B}_{i-1}\|_F^2 + \tau_i\|\mathbf{C}_{i-1}\|_F^2 + \tau_i(1 - \tau_i)\left(\|\mathbf{B}_{i-1}\|_F^2 + \|\mathbf{C}_{i-1}\|_F^2\right)\right) \\
& \leq 2\left(\|\mathbf{H}\|_F^2 + (1 + \tau_i(1 - \tau_i))\max_{\mathcal{D}_M}\|\cdot\|_F^2\right) \\
& \leq 2\left(\|\mathbf{H}\|_F^2 + \frac{5}{16}M^2d_0^2d_1^2\right),
\end{aligned}$$

observing that $\tau_i \in (0, 1]$ hence $\tau_i(1 - \tau_i) \leq \frac{1}{4}$. Thus, for $\mathbf{H} = \mathbf{B}^*$, a global minimizer of Φ ,

$$\begin{aligned}
\frac{\Phi(\mathbf{B}_k) - \Phi(\mathbf{B}^*)}{A_k} + \sum_{i=1}^k \frac{1 - L\beta_i}{2A_i\beta_i} \|\mathbf{B}_i - \mathbf{A}_i\|_F^2 & \leq \frac{\|\mathbf{C}_0 - \mathbf{B}^*\|_F^2}{2\gamma_1} \\
& + \sum_{i=1}^k A_i^{-1} \left(\delta' + L'\tau_i \left(\|\mathbf{B}^*\|_F^2 + \frac{5}{16}M^2d_0^2d_1^2 \right) \right).
\end{aligned}$$

By construction, $\sum_{i=1}^k A_i^{-1}L'\tau_i = \frac{L'}{A_k}$, and $\Phi(\mathbf{B}_k) - \Phi(\mathbf{B}^*) \geq 0$. It follows that

$$\begin{aligned}
& \min_{i=1}^k \|\beta_i^{-1}(\mathbf{B}_i - \mathbf{A}_i)\|_F^2 \\
& \leq 2 \left(\sum_{i=1}^k \frac{\beta_i(1 - L\beta_i)}{A_i} \right)^{-1} \left(\frac{\|\mathbf{C}_0 - \mathbf{B}^*\|_F^2}{2\gamma_1} + \sum_{i=1}^k A_i^{-1}\delta' + \frac{L'}{A_k} \left(\|\mathbf{B}^*\|_F^2 + \frac{5}{16}M^2d_0^2d_1^2 \right) \right).
\end{aligned}$$

As $\beta_i = \frac{L}{2}$, $\gamma_1 = \frac{1}{4L}$, and $A_i = \frac{2}{i(i+1)}$, $\sum_{i=1}^k \frac{\beta_i(1 - L\beta_i)}{A_i} = \frac{1}{4L} \sum_{i=1}^k A_i^{-1} = \frac{k(k+1)(k+2)}{24L}$, so

$$\min_{i=1}^k \|\beta_i^{-1}(\mathbf{B}_i - \mathbf{A}_i)\|_F^2 \leq \frac{96L^2}{k(k+1)(k+2)} \|\mathbf{C}_0 - \mathbf{B}^*\|_F^2 + 8L\delta' + \frac{24LL'}{N} \left(\|\mathbf{B}^*\|_F^2 + \frac{5M^2d_0^2d_1^2}{16} \right).$$

This proves the claimed result in the non-convex setting.

In the convex regime, recall from the prior discussion that we may set $L' = 0$ in the previous display, proving the claim.

E Additional Results

E.1 Proof of Lemma 1

The proof of Lemma 1 follows from the following lemma coupled with the chain rule for Fréchet differentiable maps.

Lemma 10. *Let $\mu_i \in \mathcal{P}(\mathbb{R}^{d_i})$, for $i = 0, 1$, be compactly supported with $\text{spt}(\mu_i) = S_i$. Then, the map $f \in \mathcal{C}(S_0 \times S_1) \mapsto \left(\int e^{f(\cdot, y)} d\mu_1(y), \int e^{f(x, \cdot)} d\mu_0(x) \right) \in \mathcal{C}(S_0) \times \mathcal{C}(S_1)$ is smooth with first derivative at $f \in \mathcal{C}(S_0 \times S_1)$ given by*

$$h \in \mathcal{C}(S_0 \times S_1) \mapsto \left(\int h(\cdot, y) e^{f(\cdot, y)} d\mu_1(y), \int h(x, \cdot) e^{f(x, \cdot)} d\mu_0(x) \right) \in \mathcal{C}(S_0) \times \mathcal{C}(S_1).$$

Proof. First, we show that the map $f \in \mathcal{C}(S_0 \times S_1) \mapsto e^f \in \mathcal{C}(S_0 \times S_1)$ is Fréchet differentiable with $D(e^{(\cdot)})_{[f]}(h) = he^f$. Fix $f \in \mathcal{C}(S_0 \times S_1)$ and consider

$$\lim_{\substack{h \in \mathcal{C}(S_0 \times S_1) \\ \|h\|_{\infty, S_0 \times S_1} \rightarrow 0}} \frac{\|e^{f+h} - e^f - he^f\|_{\infty, S_0 \times S_1}}{\|h\|_{\infty, S_0 \times S_1}} \leq \|e^f\|_{\infty, S_0 \times S_1} \lim_{\substack{h \in \mathcal{C}(S_0 \times S_1) \\ \|h\|_{\infty, S_0 \times S_1} \rightarrow 0}} \frac{\|e^h - 1 - h\|_{\infty, S_0 \times S_1}}{\|h\|_{\infty, S_0 \times S_1}}.$$

Fix arbitrary $(x, y) \in S_0 \times S_1$. By a Taylor expansion,

$$e^{h(x,y)} = 1 + h(x, y) + \frac{1}{2}e^{\xi(x,y)}h^2(x, y),$$

where $|\xi(x, y)| \in [0, |h(x, y)|]$ i.e. $\|\xi\|_{\infty, S_0 \times S_1} \leq \|h\|_{\infty, S_0 \times S_1}$. That is,

$$\lim_{\substack{h \in \mathcal{C}(S_0 \times S_1) \\ \|h\|_{\infty, S_0 \times S_1} \rightarrow 0}} \frac{\|e^h - 1 - h\|_{\infty, S_0 \times S_1}}{\|h\|_{\infty, S_0 \times S_1}} \leq \lim_{\|h\|_{\infty, S_0 \times S_1} \rightarrow 0} \frac{1}{2}e^{\|h\|_{\infty, S_0 \times S_1}} \|h\|_{\infty, S_0 \times S_1} = 0.$$

On the other hand, the derivative of $f \in \mathcal{C}(S_0 \times S_1) \mapsto \int f(x, y)d\mu_1(y) \in \mathcal{C}(S_0)$ at any point is given by $h \in \mathcal{C}(S_0 \times S_1) \mapsto \int h(x, y)d\mu_1(y) \in \mathcal{C}(S_0)$. The claimed expression for the first derivative then follows by the chain rule. The derivatives of this map can be computed to arbitrary order inductively by the prior argument. \square

Proof of Lemma 1. Observe that the map $(\mathbf{A}, \varphi_0, \varphi_1) \in \mathbb{R}^{d_0 \times d_1} \times \mathfrak{E} \mapsto \varphi_0 \oplus \varphi_1 - c_{\mathbf{A}} \in \mathcal{C}(S_0 \times S_1)$ is smooth with first derivative at $(\mathbf{A}, \varphi_0, \varphi_1) \in \mathbb{R}^{d_0 \times d_1} \times \mathfrak{E}$ given by

$$(\mathbf{B}, h_0, h_1) \in \mathbb{R}^{d_0 \times d_1} \times \mathfrak{E} \mapsto h_0 \oplus h_1 + 32x^T \mathbf{B}y \in \mathcal{C}(S_0 \times S_1).$$

The result then follows from Lemma 10 by applying the chain rule. \square

E.2 Compactness of \mathcal{L}

Lemma 11 (Example 2 in [42]). *Let $\varepsilon > 0$, $\mu_0 \in \mathcal{P}(\mathbb{R}^{d_0})$, $\mu_1 \in \mathcal{P}(\mathbb{R}^{d_1})$, and $\mathbf{A} \in \mathbb{R}^{d_0 \times d_1}$ be arbitrary and let $(\varphi_0^{\mathbf{A}}, \varphi_1^{\mathbf{A}})$ be EOT potentials for $\text{OT}_{\mathbf{A}, \varepsilon}(\mu_0, \mu_1)$. Then, the map $\mathcal{L} : L^2(\mu_0) \times L^2(\mu_1) \mapsto L^2(\mu_0) \times L^2(\mu_1)$ defined by*

$$\mathcal{L}(f_0, f_1) = \left(\int f_1(y) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x,y)}{\varepsilon}} d\mu_1(y), \int f_0(x) e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x,y)}{\varepsilon}} d\mu_0(x) \right),$$

is compact.

Proof. For simplicity, we prove only that

$$\mathcal{L}_2 : f \in L^2(\mu_0) \mapsto \int f(x) \xi(x, \cdot) d\mu_0(x) \in L^2(\mu_1),$$

is a compact operator for $\xi : (x, y) \in \mathbb{R}^{d_0} \times \mathbb{R}^{d_1} \mapsto e^{\frac{\varphi_0^{\mathbf{A}}(x) + \varphi_1^{\mathbf{A}}(y) - c_{\mathbf{A}}(x,y)}{\varepsilon}}$. For any $y \in \mathbb{R}^{d_1}$ and $f \in L^2(\mu_0)$, $|\mathcal{L}_2(f)(y)|^2 \leq \|f\|_{L^2(\mu_0)}^2 \int |\xi(\cdot, y)|^2 d\mu_0$, as $\xi(\cdot, y)$ is bounded on $\text{spt}(\mu_0)$ such that this operator is well-defined.

Let f_n be a bounded sequence in $L^2(\mu_0)$. By the Eberlein-Šmulian theorem [42, p. 141], up to passing to a subsequence, f_n converges weakly to $f \in L^2(\mu_0)$. For fixed $y \in \mathbb{R}^{d_1}$, $\xi(\cdot, y) \in L^2(\mu_0)$, hence $\mathcal{L}_2(f_n)(y) \rightarrow \mathcal{L}_2(f)(y)$ and it follows from the dominated convergence theorem that, for any $g \in L^2(\mu_1)$, $\int \mathcal{L}_2(f_n)g d\mu_1 \rightarrow \int \mathcal{L}_2(f)g d\mu_1$ such that $\mathcal{L}_2(f_n) \rightarrow \mathcal{L}_2(f)$ weakly in $L^2(\mu_1)$. Also, by dominated convergence,

$$\|\mathcal{L}_2(f_n)\|_{L^2(\mu_1)}^2 = \int \mathcal{L}_2(f_n)^2 d\mu_1 \rightarrow \int \mathcal{L}_2(f)^2 d\mu_1 = \|\mathcal{L}_2(f)\|_{L^2(\mu_1)}^2,$$

such that $\mathcal{L}_2(f_n) \rightarrow \mathcal{L}_2(f)$ strongly in $L^2(\mu_1)$. As f_n was an arbitrary bounded sequence in $L^2(\mu_0)$ and $\mathcal{L}_2(f_n) \rightarrow \mathcal{L}_2(f)$ strongly in $L^2(\mu_1)$ up to a subsequence, \mathcal{L}_2 is a compact operator. \square

F Blown-up figures

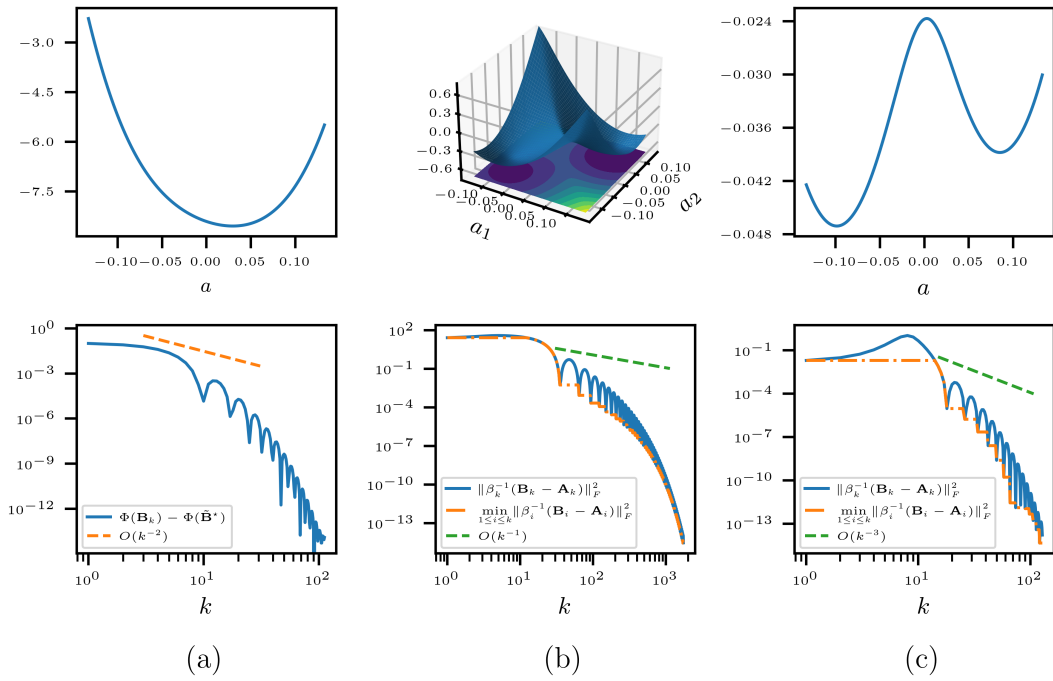


Figure 1: The top row compiles plots of Φ for the different examples described in the text. The bottom row consists of plots tracking the progress of the iterates. In (b) and (c), Algorithm 2 is initialized at $\mathbf{C}_0 = (1, 1) \times 10^{-5}$ and $\mathbf{C}_0 = 1 \times 10^{-5}$, respectively.

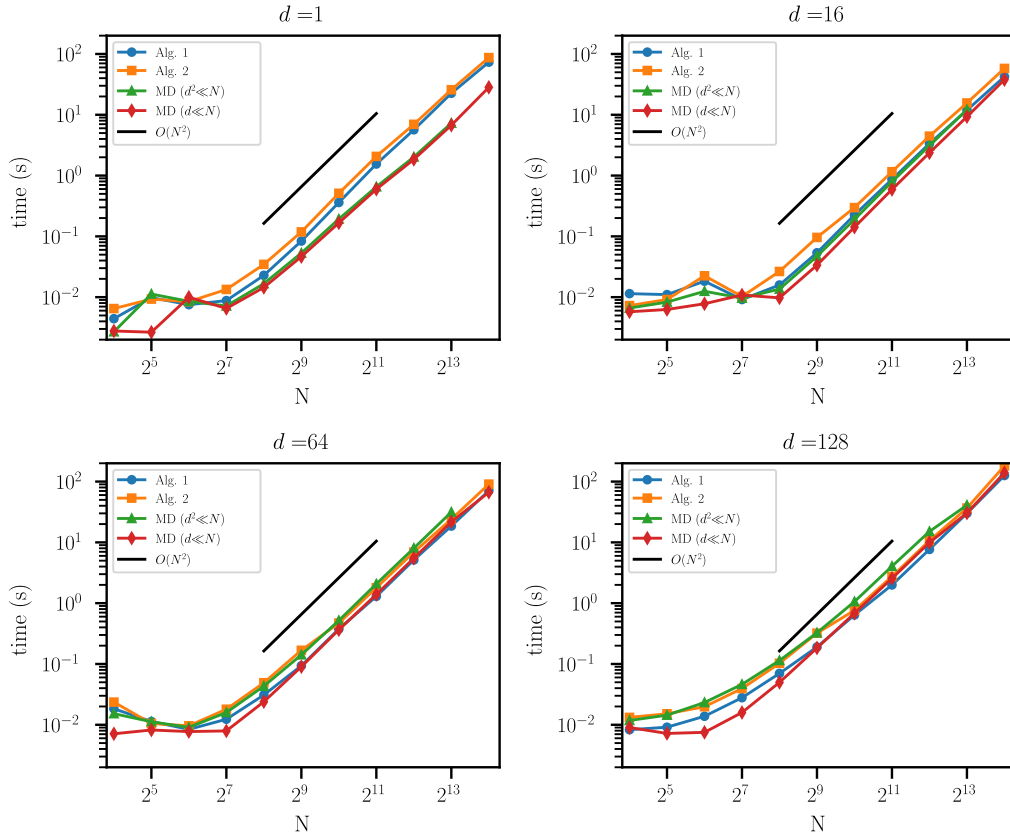


Figure 2: The various plots compile the average runtime of Algorithms 1 and 2, and two versions of the mirror descent algorithm in the convex regime for different combinations of d and N .

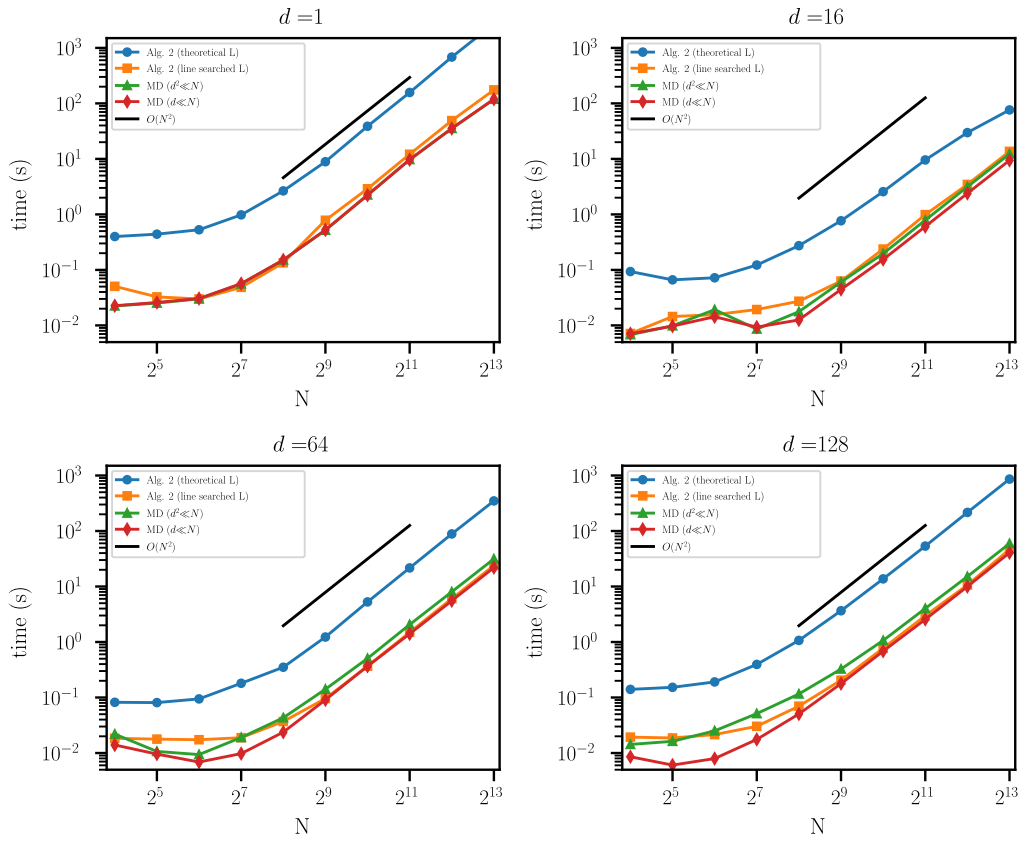


Figure 3: The various plots compile the average runtime of Algorithm 2 with the two methods for choosing L , and two versions of the mirror descent algorithm in the non-convex regime for different combinations of d and N .