

---

# Optimal Transport with Adaptive Regularisation

---

Hugues Van Assel<sup>1\*</sup>, Titouan Vayer<sup>2</sup>, Rémi Flamary<sup>3</sup>, Nicolas Courty<sup>4</sup>

<sup>1</sup>UMPA ENS Lyon, <sup>2</sup>LIP ENS Lyon, <sup>3</sup>CMAP École polytechnique,

<sup>4</sup>IRISA Université Bretagne Sud

## Abstract

Regularising the primal formulation of optimal transport (OT) with a strictly convex term leads to enhanced numerical complexity and a denser transport plan. Many formulations impose a global constraint on the transport plan, for instance by relying on entropic regularisation. As it is more expensive to diffuse mass for outlier points compared to central ones, this typically results in a significant imbalance in the way mass is spread across the points. This can be detrimental for some applications where a minimum of smoothing is required per point. To remedy this, we introduce OT with Adaptive Regularisation (OTARI), a new formulation of OT that imposes constraints on the mass going in or/and out of each point. We then showcase the benefits of this approach for domain adaptation.

## 1 Introduction

Optimal transport (OT) is a well-established framework to compare probability distributions with numerous applications in machine learning [1, 20, 21]. Discrete OT seeks a transportation plan that minimizes the total transportation cost between samples from the source and target distributions. In the absence of regularisation, this optimal OT plan is inherently sparse. Regularising OT with a strictly convex term is a widely adopted practical approach, leading to reduced numerical complexity and more diffuse OT plans [21]. As an illustration, the prominent entropic regularisation [9] leads to a dense plan. In some applications, the smoothing effect induced by the regularisation has a primary importance on its own. A key example is the construction of doubly stochastic affinity matrices for clustering and dimensionality reduction [15, 27], where smoothing enables connecting to neighbor points. Another is domain adaptation [8] where smoothed OT often results in enhanced performance when compared to non-regularised ones (see for instance Table 1). Many OT regularisation schemes on the primal formulation impose a constraint on the overall transport plan. Consequently, central data points tend to exhibit a denser (or more diffuse) transport plan compared to extreme (or outlier) data points, for which diffusion is more costly. As a result, the latter points receive limited benefits from the smoothing effect introduced by the regularizer as shown in the left side of Figure 1. This partly explains OT’s significant sensitivity to outliers in many applications [19, 7]. To remedy this, one needs to constrain the transport plan in a pointwise manner. Note that this has recently been explored for constructing affinity matrices [25] (*i.e.* symmetric OT setting) leading to enhanced noise robustness and clustering abilities.

**Contributions.** In this work, we develop a new formulation of OT, called OT with Adaptive Regularisation (OTARI), allowing to control, for any strictly convex function  $\psi$ , the value of  $\psi$  on each row and/or column of the OT plan. We then show the advantages of OTARI over usual regularised OT on domain adaptation tasks, focusing particularly on the negative entropy and the  $\ell_2$  norm respectively associated with entropic [9] and quadratic [4] optimal transport.

---

\*hugues.van\_assel@ens-lyon.fr

## 2 Regularised Optimal Transport

We first introduce the discrete OT problem before presenting regularised formulations and associated algorithms. Let  $X_S = \{\mathbf{x}_i^S \in \mathbb{R}^d\}_{i=1}^{N_S}$  and  $X_T = \{\mathbf{x}_i^T \in \mathbb{R}^d\}_{i=1}^{N_T}$  denote the sets of respectively source and target point locations. The discrete *Monge-Kantorovitch* problem [13] focuses on the optimal allocation strategy to transport the empirical measure  $\mu_S = \frac{1}{N_S} \sum_{i=1}^{N_S} a_i \delta_{\mathbf{x}_i^S}$  onto  $\mu_T = \frac{1}{N_T} \sum_{i=1}^{N_T} b_i \delta_{\mathbf{x}_i^T}$  where  $\mathbf{a} \in \Delta^{N_S}$  and  $\mathbf{b} \in \Delta^{N_T}$ . It consists in computing a *coupling*  $\mathbf{P} \in \mathbb{R}_+^{N_S \times N_T}$  i.e. a joint probability measure over  $X_S \times X_T$  solving the linear program

$$\min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle \quad (\text{OT})$$

where  $\Pi(\mathbf{a}, \mathbf{b})$  is the transport polytope with marginals  $(\mathbf{a}, \mathbf{b})$  and the *cost matrix*  $\mathbf{C} \in \mathbb{R}_+^{N_S \times N_T}$  encodes the pairwise distances between the source and target samples. One can typically consider the squared Euclidean distance  $C_{ij} = \|\mathbf{x}_i^S - \mathbf{x}_j^T\|_2^2$  or any Riemannian distance over a manifold [26].

To enable faster algorithmic resolution as well as smoother solutions, one can rely on a strictly convex regulariser  $\psi : \mathbb{R}^{N_S} \rightarrow \mathbb{R}$ . It amounts to solving  $\min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon^* \sum_i \psi(\mathbf{P}_{i\cdot})$  where  $\varepsilon^* > 0$ . Interestingly, regularised OT can also be framed using a convex constraint as follows

$$\min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle \quad \text{s.t.} \quad \mathbf{P} \in \bar{\mathcal{B}}(\eta) \quad (\text{ROT})$$

where  $\bar{\mathcal{B}}(\eta) := \{\mathbf{P} \text{ s.t. } \sum_i \psi(\mathbf{P}_{i\cdot}) \leq \eta\}$ . Note that the previously introduced  $\varepsilon^*$  is the optimal dual variable associated with the constraint  $\bar{\mathcal{B}}(\eta)$  in the above equivalent formulation. Throughout, we make the following assumption on  $\psi$ .

**Assumption 1.** Let  $\psi : \text{dom}(\psi) \rightarrow \mathbb{R} \cup \{\infty\}$  be strictly convex and differentiable on the interior of its domain  $\text{dom}(\psi) \subset \mathbb{R}_+^{N_S}$ .

In what follows, we denote by  $\psi(\mathbf{P}) = (\psi(\mathbf{P}_{1\cdot}), \dots, \psi(\mathbf{P}_{N_S\cdot}))^\top$ . We introduce  $\psi^* := \mathbf{p} \rightarrow \sup_{\mathbf{q} \in \text{dom}(\psi)} \langle \mathbf{p}, \mathbf{q} \rangle - \psi(\mathbf{q})$  the convex conjugate of  $\psi$  [23]. Note that when  $\psi$  is strictly convex, this supremum is uniquely achieved and from Danskin's theorem [10]:  $\nabla \psi^*(\mathbf{p}) = \arg \max_{\mathbf{q} \in \text{dom}(\psi)} \langle \mathbf{p}, \mathbf{q} \rangle - \psi(\mathbf{q})$ . We show in Appendix B.1 that when  $\varepsilon^* > 0$ , i.e. when the constraint  $\mathbf{P} \in \bar{\mathcal{B}}(\eta)$  is active, (ROT) is solved for  $\mathbf{P}^* = \nabla \psi^*((\mathbf{C} - \boldsymbol{\lambda}^* \oplus \boldsymbol{\mu}^*)/\varepsilon^*)^1$  where  $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*, \varepsilon^*)$  solve the following dual problem

$$\max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \varepsilon > 0} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle + \varepsilon \left( \sum_i \psi^*((\mathbf{C}_{i\cdot} - \lambda_i \mathbf{1} - \boldsymbol{\mu})/\varepsilon) - \eta \right). \quad (\text{Dual-ROT})$$

The above objective is concave thus the problem can be solved exactly using e.g. BFGS [17] or ADAM [14]. As a complementary view, one can also frame (ROT) as a  $\psi$ -Bregman projection over a convex set. The  $\psi$ -Bregman divergence is defined as  $D_\psi(\mathbf{P}, \mathbf{Q}) := \psi(\mathbf{P}) - \psi(\mathbf{Q}) - \langle \mathbf{P} - \mathbf{Q}, \nabla \psi(\mathbf{Q}) \rangle$ . The solution of (ROT) can then be expressed as  $\mathbf{P}^* = \text{Proj}_{\Pi(\mathbf{a}, \mathbf{b}) \cap \bar{\mathcal{B}}(\eta)}^{D_\psi}(\mathbf{K}_\sigma)$  where  $\mathbf{K}_\sigma := \nabla \psi^*(-\mathbf{C}/\sigma)$  for any  $\sigma < \varepsilon^*$  (see Appendix B.1 for details). The key benefit of the above result is that it enables solving (ROT) with alternating Bregman projection schemes [3].

In this work, we focus specifically on certain Bregman divergences: the Kullback Leibler (KL) divergence and the squared Euclidean distance. The first reads  $D_{\text{KL}}(\mathbf{P}|\mathbf{Q}) = \langle \mathbf{P}, \log(\mathbf{P} \oslash \mathbf{Q}) - \mathbf{1}\mathbf{1}^\top \rangle$  with associated negative entropy  $\psi_{\text{KL}}(\mathbf{p}) = \langle \mathbf{p}, \log \mathbf{p} - \mathbf{1} \rangle$ . In this case, (ROT) boils down to entropic OT and solved for  $\text{Proj}_{\Pi(\mathbf{a}, \mathbf{b})}^{\text{KL}}(\mathbf{K}_{\varepsilon^*})$  where  $\mathbf{K}_{\varepsilon^*} = \nabla \psi_{\text{KL}}^*(-\mathbf{C}/\varepsilon^*) = \exp(-\mathbf{C}/\varepsilon^*)$  is a Gibbs kernel. This projection is well-known as the *static Schrödinger bridge* [16] referring to statistical physics where it first appeared [24], and can be computed efficiently using the Sinkhorn algorithm [9]. For the squared Euclidean distance, we define  $\psi_2(\mathbf{p}) = \frac{1}{2} \|\mathbf{p}\|_2^2$ . The associated problem (1) is usually referred to as quadratic OT [18] and can yield sparse OT plans unlike entropic.

## 3 Optimal Transport with Adaptive Regularisation

In this section, we present a new formulation of OT that imposes constraints on each row of the OT plan. We begin by introducing the set of matrices with *point-wise* constraints. To set the upper bound,

<sup>1</sup>We use the notation  $\nabla \psi^*(\mathbf{P}) := (\nabla \psi^*(\mathbf{P}_{1\cdot}), \dots, \nabla \psi^*(\mathbf{P}_{N_S\cdot}))^\top$ .

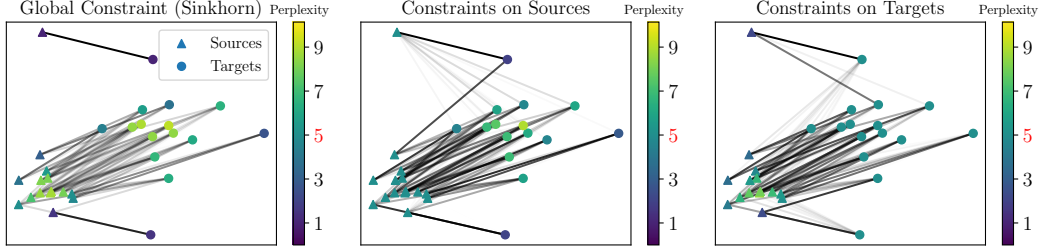


Figure 1: Entropic OT plans ( $\xi = 5$ ) with global constraint, pointwise constraints on sources and then on targets. The three plans have the same global entropy. The color of each source (resp. target) point shows the perplexity (exponential of entropy) of the associated row (resp. column) of the OT plan.

we rely on the perplexity parameter  $\xi$  [25] that can be interpreted as the number of effective neighbors for each point. Concretely, we define  $\mathbf{e}_\xi = \frac{1}{\xi} (\mathbb{1}_{i \leq \xi})_i$  and

$$\mathcal{B}_\psi(\xi) := \{\mathbf{P} \geq \mathbf{0} \text{ s.t. } \forall i, \psi(\mathbf{P}_{i,:}) \leq \psi(\mathbf{e}_\xi)\}. \quad (1)$$

Note that  $\psi_{\text{KL}}(\mathbf{p}_\xi) = -(\log \xi + 1)$  and  $\psi_2(\mathbf{p}_\xi) = 1/\xi$ . We now define Optimal Transport with Adaptive Regularisation (OTARI) as the generalization of (ROT) to the case where the strictly convex constraint is given by  $\mathcal{B}_\psi(\xi)$ . Similarly to Proposition 1 (Appendix B.1), we can frame OTARI as a  $\psi$ -Bregman projection of  $\mathbf{K}_\sigma = \nabla \psi^*(-\mathbf{C}/\sigma)$  or solve it using dual ascent.

**Proposition 2.** *Let  $(\mathbf{a}, \mathbf{b}, \xi)$  be such that  $\Pi(\mathbf{a}, \mathbf{b}) \cap \mathcal{B}_\psi(\xi)$  has an interior point and let  $\mathbf{P}^*$  solve*

$$\min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle \quad \text{s.t.} \quad \mathbf{P} \in \mathcal{B}_\psi(\xi). \quad (\text{OTARI-s})$$

*Let  $\boldsymbol{\varepsilon}^*$  be the optimal dual variable associated with the constraint  $\mathbf{P} \in \mathcal{B}_\psi(\xi)$ . If  $\boldsymbol{\varepsilon}^* > \mathbf{0}$ , then it holds  $\mathbf{P}^* = \text{Proj}_{\Pi(\mathbf{a}, \mathbf{b}) \cap \mathcal{B}_\psi(\xi)}^{D_\psi}(\mathbf{K}_\sigma)$  for any  $0 < \sigma \leq \min_i \varepsilon_i^*$ . Moreover it holds  $\mathbf{P}^* = \nabla \psi^*(\text{diag}(\boldsymbol{\varepsilon}^*)^{-1}(\mathbf{C} - \boldsymbol{\lambda}^* \oplus \boldsymbol{\mu}^*))$  where  $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*, \boldsymbol{\varepsilon}^*)$  solve the following dual*

$$\max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\varepsilon} > \mathbf{0}} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle + \langle \boldsymbol{\varepsilon}, \psi^*(\text{diag}(\boldsymbol{\varepsilon})^{-1}(\mathbf{C} - \boldsymbol{\lambda} \oplus \boldsymbol{\mu})) - \psi(\mathbf{e}_\xi) \mathbf{1} \rangle. \quad (\text{Dual-OTARI-s})$$

---

**Algorithm 1** *Dijkstra* for solving (OTARI-d)

---

- 1: **Input:**  $\mathbf{C}, \psi(\cdot), \xi^a, \xi^b, \varepsilon, \mathbf{a}, \mathbf{b}$
  - 2:  $(\mathbf{P}_b, \boldsymbol{\Xi}, \boldsymbol{\Theta}) \leftarrow (\nabla \psi^*(-\mathbf{C}/\varepsilon), \mathbf{0}, \mathbf{0})$
  - 3: **while** not converged **do**
  - 4:  $\mathbf{P}_a \leftarrow \text{Proj}_{\Pi(\mathbf{a})}^{D_\psi}(\mathbf{P}_b)$
  - 5:  $\bar{\mathbf{P}}_a \leftarrow \text{Proj}_{\mathcal{B}_\psi(\xi^a)}^{D_\psi} \circ \nabla \psi^*(\nabla \psi(\mathbf{P}_a) + \boldsymbol{\Xi})$
  - 6:  $\boldsymbol{\Xi} \leftarrow \boldsymbol{\Xi} + \nabla \psi(\mathbf{P}_a) - \nabla \psi(\bar{\mathbf{P}}_a)$
  - 7:  $\mathbf{P}_b^\top \leftarrow \text{Proj}_{\Pi(\mathbf{b})}^{D_\psi}(\bar{\mathbf{P}}_a^\top)$
  - 8:  $\bar{\mathbf{P}}_b^\top \leftarrow \text{Proj}_{\mathcal{B}_\psi(\xi^b)}^{D_\psi} \circ \nabla \psi^*((\nabla \psi(\mathbf{P}_b) + \boldsymbol{\Theta})^\top)$
  - 9:  $\boldsymbol{\Theta} \leftarrow \boldsymbol{\Theta} + \nabla \psi(\mathbf{P}_b) - \nabla \psi(\bar{\mathbf{P}}_b)$
  - 10: **end while**
  - 11: **Output:**  $\bar{\mathbf{P}}_b$
- 

(OTARI-d) that consists of projecting  $\mathbf{K}_\sigma$  onto the nonempty set  $\mathcal{B}_\psi(\xi^a) \cap \mathcal{B}_\psi^\top(\xi^b)$  where we defined  $\mathcal{B}_\psi^\top(\xi) = \{\mathbf{P}^\top \in \mathcal{B}_\psi(\xi)\}$  thus ensuring sufficient smoothing for both rows and columns. Such projection can be computed using alternating Bregman projections, whose convergence has been well-studied [6, 3]. As we generally do not have access to a closed form for the projection onto the transport polytope  $\Pi(\mathbf{a}, \mathbf{b})$ , it is common to alternate projection onto  $\Pi(\mathbf{a})$  and  $\Pi(\mathbf{b})$  separately (see *e.g.* the seminal Sinkhorn algorithm [9]). We extend this scheme by adding projection steps into the pointwise constraints  $\mathcal{B}_\psi(\xi)$  for both  $\mathbf{P}$  and  $\mathbf{P}^\top$ . As this set is not affine, one needs to resort to the Dykstra procedure [11] that can be applied to a broad class of Bregman divergences [2], as shown in Algorithm 1. In Appendix C, we provide the form of the projections for  $\psi_{\text{KL}}$  and  $\psi_2$ . A key benefit of decoupling both row and column constraints is that projection onto  $\mathcal{B}_\psi(\xi)$  exhibits a simple structure where the rows can be decoupled into independent subproblems.

According to Proposition 2, one can solve (OTARI-s) using either alternating projections or dual ascent. When  $\boldsymbol{\varepsilon}^* > \mathbf{0}$ , meaning that all constraints are active *i.e.*  $\forall i, \psi(\mathbf{P}_{i,:}^*) = \psi(\mathbf{p}_\xi)$ , dual ascent is usually faster. However, if  $\boldsymbol{\varepsilon}^*$  has null components, one can still rely on  $\text{Proj}_{\Pi(\mathbf{a}, \mathbf{b}) \cap \mathcal{B}_\psi(\xi)}^{D_\psi}(\mathbf{K}_\varepsilon)$  to provide an approximate solution as alternating Bregman projections are always guaranteed to converge.

Note that we can impose the pointwise constraint equivalently on the rows or the columns of the OT plan. Hence (OTARI-t) can be defined by imposing the constraint on the target samples *i.e.*  $\mathbf{P}^\top \in \mathcal{B}_\psi(\xi)$ . We also propose a doubly constrained formulation called

## 4 Application to Domain Adaptation

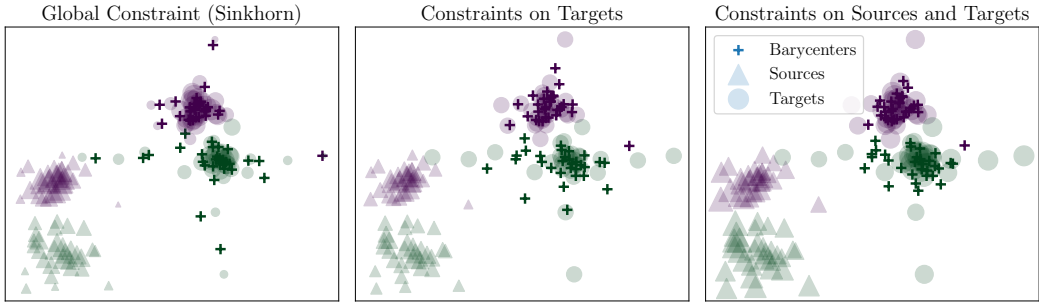


Figure 2: Toy domain adaptation scenario with entropic OT plans ( $\xi = 10$ ) with various constraints. The size of the point is proportional to the associated entropy. When using Sinkhorn, barycentric mapping match outlier points since the OT plan is less diffuse for these points. In turn, using pointwise constraints concentrate the mapped points in high-density regions, thus giving more robust estimates for the mappings onto the target domain.

	OT	EOT	EOTARI-s	EOTARI-t	EOTARI-d
MNIST $\rightarrow$ USPS ( $\xi = 30$ )	53.1(5.4)	64.2(2.8)	65.0(5.3)	66.4(3.5)	<b>67.4(2.9)</b>
MNIST $\rightarrow$ USPS ( $\xi = 300$ )	53.1(5.4)	68.8(3.1)	70.8(4.2)	70.2(3.4)	<b>72.6(5.1)</b>
USPS $\rightarrow$ MNIST ( $\xi = 30$ )	59.1(4.9)	60.8(5.4)	61.6(4.4)	<b>62.6(3.0)</b>	61.0(4.7)
USPS $\rightarrow$ MNIST ( $\xi = 300$ )	59.1(4.9)	59.8(1.6)	61.0(2.3)	<b>61.6(3.0)</b>	58.8(2.3)
	OT	QOT	QOTARI-s	QOTARI-t	QOTARI-d
MNIST $\rightarrow$ USPS ( $\xi = 30$ )	53.1(5.4)	68.3(3.9)	68.3(3.6)	<b>69.3(4.7)</b>	68.1(4.6)
MNIST $\rightarrow$ USPS ( $\xi = 300$ )	53.1(5.4)	60.7(1.5)	<b>67.0(2.4)</b>	65.5(2.3)	65.8(2.5)
USPS $\rightarrow$ MNIST ( $\xi = 30$ )	59.1(4.9)	60.4(3.5)	<b>62.8(3.7)</b>	59.6(2.7)	61.6(3.1)
USPS $\rightarrow$ MNIST ( $\xi = 300$ )	59.1(4.9)	59.2(3.4)	60.1(3.0)	<b>62.0(3.7)</b>	61.5(3.8)

Table 1: Domain adaptation 1-NN classification scores for OT (unregularised), EOT (entropic), EOTARI (entropic OTARI), QOT (quadratic), QOTARI (quadratic OTARI) for  $\xi = 30$  and  $\xi = 300$ .

In this section, we evaluate OTARI on a domain adaptation task where the goal is to transport labeled data points to a target domain where a classifier is trained. Mapping onto the target domain is performed through a barycentric mapping of the form: for any  $i \in \llbracket N_S \rrbracket$ ,  $\hat{\mathbf{x}}_i = \frac{1}{a_i} \sum_j T_{ij} \mathbf{x}_j^T$ . Looking at Figure 2, one can notice that using OTARI for domain adaptation yields a mapping that is concentrated in high-density (thus more faithful) regions of the target domain. On the opposite, when using globally constrained OT (left side of Figure 2), the barycentric mapping associated with an outlier is concentrated on the outlier’s position. For the experiments, we take  $\mathbf{C}$  as the squared Euclidean distance computed from raw images of the handwritten digit classification benchmark MNIST-USPS. Following the standard practice in OT-based domain adaptation [12], we map the source to the target samples and then train a 1-NN classifier on the barycentric mappings with source labels. We compute the outcomes across 10 independent trials. In each of these experiments, the target data is partitioned into a 90% training and 10% testing split, with OT barycentric mappings and 1-NN classifiers exclusively applied to the training set. Mean scores and standard deviations are displayed in Table 1. The latter shows that adaptive regularisation consistently outperforms global regularisation (set such that the average perplexity is  $\xi$  for a fair comparison) with significant performance gains in some settings (see *e.g.* MNIST  $\rightarrow$  USPS ( $\xi = 300$ ) with the quadratic regularisation).

## 5 Conclusion

In this work, we presented a versatile framework to control the value of any OT regulariser in source or/and target locations. We showed encouraging preliminary results for domain adaptation that will be investigated in upcoming works. One could also extend OTARI to continuous distributions and apply it to OT mapping estimation [22]. Other interesting directions include investigating optimization algorithms that can avoid quadratic memory complexity.

## References

- [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- [2] Heinz H Bauschke and Adrian S Lewis. Dykstras algorithm with bregman projections: A convergence proof. *Optimization*, 48(4):409–427, 2000.
- [3] Jean-David Benamou, Guillaume Carlier, Marco Cuturi, Luca Nenna, and Gabriel Peyré. Iterative bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138, 2015.
- [4] Mathieu Blondel, Vivien Seguy, and Antoine Rolet. Smooth and sparse optimal transport. In *International conference on artificial intelligence and statistics*, pages 880–889. PMLR, 2018.
- [5] Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [6] Yair Censor and Simeon Reich. The dykstra algorithm with bregman projections. *Communications in Applied Analysis*, 2(3):407–420, 1998.
- [7] Ching-Yao Chuang, Stefanie Jegelka, and David Alvarez-Melis. InfoOT: Information maximizing optimal transport. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 6228–6242. PMLR, 23–29 Jul 2023.
- [8] Nicolas Courty, Rémi Flamary, Amaury Habrard, and Alain Rakotomamonjy. Joint distribution optimal transportation for domain adaptation. *Advances in neural information processing systems*, 30, 2017.
- [9] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Neural Information Processing Systems (NeurIPS)*, 26, 2013.
- [10] John M Danskin. The theory of max-min, with applications. *SIAM Journal on Applied Mathematics*, 14(4):641–664, 1966.
- [11] Richard L Dykstra. An algorithm for restricted least squares regression. *Journal of the American Statistical Association*, 78(384):837–842, 1983.
- [12] Rémi Flamary, Cédric Févotte, Nicolas Courty, and Valentin Emiya. Optimal spectral transportation with application to music transcription. *Neural Information Processing Systems (NeurIPS)*, 29, 2016.
- [13] L Kantorovich. On the transfer of masses (in russian). In *Doklady Akademii Nauk*, volume 37, page 227, 1942.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [15] Boris Landa, Ronald R Coifman, and Yuval Kluger. Doubly stochastic normalization of the gaussian kernel is robust to heteroskedastic noise. *SIAM journal on mathematics of data science*, 3(1):388–413, 2021.
- [16] Christian Léonard. A survey of the schrödinger problem and some of its connections with optimal transport. *arXiv preprint arXiv:1308.0215*, 2013.
- [17] Dong C Liu and Jorge Nocedal. On the limited memory bfgs method for large scale optimization. *Mathematical programming*, 45(1-3):503–528, 1989.
- [18] Dirk A Lorenz, Paul Manns, and Christian Meyer. Quadratically regularized optimal transport. *Applied Mathematics & Optimization*, 83(3):1919–1949, 2021.
- [19] Debarghya Mukherjee, Aritra Guha, Justin M Solomon, Yuekai Sun, and Mikhail Yurochkin. Outlier-robust optimal transport. In *International Conference on Machine Learning*, pages 7850–7860. PMLR, 2021.

- [20] Sherjil Ozair, Corey Lynch, Yoshua Bengio, Aaron Van den Oord, Sergey Levine, and Pierre Sermanet. Wasserstein dependency measure for representation learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- [21] Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- [22] Aram-Alexandre Pooladian and Jonathan Niles-Weed. Entropic estimation of optimal transport maps. *arXiv preprint arXiv:2109.12004*, 2021.
- [23] R Tyrrell Rockafellar. *Convex analysis*, volume 11. Princeton university press, 1997.
- [24] Erwin Schrödinger. *Über die umkehrung der naturgesetze*. Verlag der Akademie der Wissenschaften in Kommission bei Walter De Gruyter u . . . , 1931.
- [25] Hugues Van Assel, Titouan Vayer, Rémi Flamary, and Nicolas Courty. Snekhorn: Dimension reduction with symmetric entropic affinities. *arXiv preprint arXiv:2305.13797*, 2023.
- [26] Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2009.
- [27] Ron Zass and Amnon Shashua. Doubly stochastic normalization for spectral clustering. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Neural Information Processing Systems (NeurIPS)*. MIT Press, 2006.

## A Notations

We adopt the conventions that  $0/0 = 0$ ,  $0 \log(0) = 0$  and  $x/0 = \infty$  for  $x > 0$ .  $\exp$ ,  $\log$  applied to vectors/matrices are taken element-wise.  $\mathbf{1}$  is the all-one vector whose size depends on the context.  $\langle \cdot, \cdot \rangle$  is the standard inner product for matrices/vectors.  $P_{ij}$  denotes the entry at position  $(i, j)$  of a matrix  $\mathbf{P}$  while  $\mathbf{P}_{i\cdot}$  and denotes the  $i$ -th row.  $\mathbf{P} \geq \mathbf{0}$  means that for any  $(i, j)$ ,  $P_{ij} \geq 0$ .  $\odot$  (*resp.*  $\oslash$ ) stands for element-wise multiplication (*resp.* division) between vectors/matrices. For  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ ,  $\mathbf{a} \oplus \mathbf{b} \in \mathbb{R}^{n \times n}$  is  $(a_i + b_j)_{ij}$ . For  $\alpha \in \mathbb{R}$ ,  $\mathbf{p}^{\odot \alpha}$  and  $\mathbf{P}^{\odot \alpha}$  denote element-wise exponentiation *i.e.*  $[\mathbf{p}^{\odot \alpha}]_i = p_i^\alpha$ .  $[\mathbf{P}]_+$  is the element-wise positive part with  $\max(0, P_{ij})$  in position  $(i, j)$ . For  $n \in \mathbb{N}$ ,  $\Delta^n$  is the probability simplex  $\{\mathbf{p} \in \mathbb{R}_+^n \text{ s.t. } \sum_i p_i = 1\}$ . For  $\mathbf{a} \in \Delta^{N_S}$  and  $\mathbf{b} \in \Delta^{N_T}$ ,  $\Pi(\mathbf{a}, \mathbf{b}) = \{\mathbf{P} \in \mathbb{R}_+^{N_S \times N_T} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{a} \text{ and } \mathbf{P}^\top \mathbf{1} = \mathbf{b}\}$  is the transport polytope with marginals  $(\mathbf{a}, \mathbf{b})$  while  $\Pi(\mathbf{a}) = \{\mathbf{P} \in \mathbb{R}_+^{N_S \times N_T} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{a}\}$  is the semi-relaxed transport polytope. For a set  $\mathcal{E}$  and a divergence  $D$ ,  $\text{Proj}_{\mathcal{E}}^D(\mathbf{K}) = \arg \min_{\mathbf{P} \in \mathcal{E}} D(\mathbf{P}|\mathbf{K})$ .

## B Proofs of the Optimal Transport Solutions

### B.1 Optimal Transport with Global Constraint

**Proposition 1.** Let  $\psi : \mathbb{R}^{N_S} \rightarrow \mathbb{R}$  satisfy Assumption 1. We define  $\bar{\mathcal{B}}(\eta) := \{\mathbf{P} \text{ s.t. } \sum_i \psi(\mathbf{P}_{i\cdot}) \leq \eta\}$ . Let  $(\mathbf{a}, \mathbf{b}, \eta)$  be such that  $\Pi(\mathbf{a}, \mathbf{b}) \cap \bar{\mathcal{B}}(\eta)$  has an interior point and let  $\mathbf{P}^*$  be a solution of

$$\min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle \quad \text{s.t.} \quad \mathbf{P} \in \bar{\mathcal{B}}(\eta). \quad (\text{ROT})$$

Let  $\varepsilon^*$  be an optimal dual variable associated with  $\mathbf{P} \in \bar{\mathcal{B}}(\eta)$ . If  $\varepsilon^* > 0$ , then for any  $0 < \sigma \leq \varepsilon^*$ , it holds  $\mathbf{P}^* = \text{Proj}_{\Pi(\mathbf{a}, \mathbf{b}) \cap \bar{\mathcal{B}}(\eta)}^{D_\psi}(\mathbf{K}_\sigma)$  where  $\mathbf{K}_\sigma := \nabla \psi^*(-\mathbf{C}/\sigma)$ . One also has  $\mathbf{P}^* = \nabla \psi^*((\mathbf{C} - \lambda^* \oplus \mu^*)/\varepsilon^*)$  where  $(\lambda^*, \mu^*, \varepsilon^*)$  solve

$$\max_{\lambda, \mu, \varepsilon > 0} \langle \lambda, \mathbf{a} \rangle + \langle \mu, \mathbf{b} \rangle + \varepsilon \left( \sum_i \psi^*((\mathbf{C}_{i\cdot} - \lambda_i \mathbf{1} - \mu)/\varepsilon) - \eta \right). \quad (\text{Dual-ROT})$$

*Proof.* We first show that  $\mathbf{P}^* = \text{Proj}_{\Pi(\mathbf{a}, \mathbf{b}) \cap \bar{\mathcal{B}}(\eta)}^{D_\psi}(\mathbf{K}_\varepsilon)$  before focusing on the dual problem.

#### Part I : Proof of the Bregman projection.

Simplifying the constant terms  $\text{Proj}_{\Pi(\mathbf{a}, \mathbf{b}) \cap \bar{\mathcal{B}}(\eta)}^{D_\psi}(\mathbf{K}_\varepsilon)$  boils down to the following problem

$$\min_{\mathbf{P}} \sum_i \psi(\mathbf{P}_{i\cdot}) - \langle \mathbf{P}, \nabla \psi(\mathbf{K}_\sigma) \rangle \quad (2)$$

$$\text{s.t.} \quad \sum_i \psi(\mathbf{P}_{i\cdot}) \leq \eta \quad (3)$$

$$\mathbf{P}\mathbf{1} = \mathbf{a}, \quad \mathbf{P}^\top \mathbf{1} = \mathbf{b} \quad (4)$$

$$\mathbf{P} \in \mathbb{R}_+^{n \times n}. \quad (5)$$

This problem is convex and strictly feasible. Strong duality holds thanks to Slater's constraint qualification. Therefore the KKT conditions [5] are necessary and sufficient conditions for optimality. The Lagrangian can be expressed as

$$\mathcal{L}(\mathbf{P}, \nu, \lambda, \mu) = \sum_i \psi(\mathbf{P}_{i\cdot}) - \langle \mathbf{P}, \nabla \psi(\mathbf{K}_\sigma) \rangle + \nu \left( \sum_i \psi(\mathbf{P}_{i\cdot}) - \eta \right) \quad (6)$$

$$+ \langle \lambda, \mathbf{a} - \mathbf{P}\mathbf{1} \rangle + \langle \mu, \mathbf{b} - \mathbf{P}^\top \mathbf{1} \rangle - \langle \Omega, \mathbf{P} \rangle. \quad (7)$$

Any optimal primal-dual variables  $(\mathbf{P}^*, \nu^*, \lambda^*, \mu^*, \Omega^*)$  satisfies

$$\nabla_{\mathbf{P}} \mathcal{L}(\mathbf{P}^*, \nu^*, \lambda^*, \mu^*) = (\nu^* + 1) \nabla \psi(\mathbf{P}^*) - \nabla \psi(\mathbf{K}_\sigma) - \lambda^* \oplus \mu^* - \Omega^* = \mathbf{0}. \quad (8)$$

Note that by definition  $\mathbf{K}_\sigma = \nabla\psi^*(-\mathbf{C}/\epsilon)$  and thus  $\nabla\psi(\mathbf{K}_\sigma) = \nabla\psi[\nabla\psi^*(-\mathbf{C}/\sigma)] = -\mathbf{C}/\sigma$  [23]. Thus we have

$$\mathbf{C} + \sigma(\nu^* + 1)\nabla\psi(\mathbf{P}^*) - \sigma\boldsymbol{\lambda}^* \oplus \boldsymbol{\mu}^* - \sigma\boldsymbol{\Omega}^* = \mathbf{0}. \quad (9)$$

Similarly, for the optimal transport problem (1)

$$\min_{\mathbf{T}} \langle \mathbf{T}, \mathbf{C} \rangle \quad \text{s.t.} \quad \mathbf{T} \in \Pi(\mathbf{a}, \mathbf{b}) \cap \bar{\mathcal{B}}(\eta), \quad (10)$$

we get the optimality condition for optimal primal-dual variables  $(\mathbf{T}^*, \varepsilon^*, \boldsymbol{\rho}^*, \boldsymbol{\kappa}^*, \boldsymbol{\Lambda}^*)$

$$\mathbf{C} + \varepsilon^*\nabla\psi(\mathbf{T}^*) - \boldsymbol{\rho}^* \oplus \boldsymbol{\kappa}^* - \boldsymbol{\Lambda}^* = \mathbf{0}. \quad (11)$$

Focusing on the original Bregman projection problem, we can then consider

$$\begin{cases} \boldsymbol{\lambda} &= \boldsymbol{\rho}^*/\sigma \\ \boldsymbol{\mu} &= \boldsymbol{\kappa}^*/\sigma \\ \nu &= \varepsilon^*/\sigma - 1 \\ \boldsymbol{\Omega} &= \boldsymbol{\Lambda}^*/\sigma \\ \mathbf{P} &= \mathbf{T}^*. \end{cases} \quad (12)$$

With the above choice,

- $(\mathbf{P}, \nu, \boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\Omega})$  satisfies the first order optimality condition.
- $\mathbf{P} = \mathbf{T}^* \in \Pi(\mathbf{a}, \mathbf{b}) \cap \mathcal{E}(\eta)$  thus the primal constraint is satisfied.
- $\sigma \leq \varepsilon^*$  implies that  $\nu \geq 0$  and  $\boldsymbol{\Omega}$  has positive entries thereby dual constraints are satisfied.
- $\varepsilon^* \neq 0$  thus by complementary slackness  $\psi(\mathbf{T}^*) = \eta$  hence complementary slackness is also verified for  $\mathbf{P}$  since  $\mathbf{P} = \mathbf{T}^*$ .

Therefore the KKT conditions are met hence  $(\mathbf{P}, \nu, \boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\Omega}) = (\mathbf{P}^*, \nu^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*, \boldsymbol{\Omega}^*)$  and  $\mathbf{P}^* = \mathbf{T}^*$ .  $\square$

## Part II : Proof of dual ascent.

The optimal dual variables  $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*, \varepsilon^*)$  solve the following problem

$$\max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \varepsilon > 0} \min_{\mathbf{P} \geq 0} \langle \mathbf{P}, \mathbf{C} \rangle + \langle \boldsymbol{\lambda}, \mathbf{a} - \mathbf{P}\mathbf{1} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} - \mathbf{P}^\top \mathbf{1} \rangle + \varepsilon \left( \sum_i \psi(\mathbf{P}_{i:}) - \eta \right) \quad (13)$$

$$= \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \varepsilon > 0} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle - \varepsilon\eta + \min_{\mathbf{P} \geq 0} \langle \mathbf{P}, \mathbf{C} - \boldsymbol{\lambda}\mathbf{1}^\top - \mathbf{1}\boldsymbol{\mu}^\top \rangle + \varepsilon \sum_i \psi(\mathbf{P}_{i:}) \quad (14)$$

$$= \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \varepsilon > 0} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle - \varepsilon\eta + \min_{\mathbf{P} \geq 0} \sum_i \langle \mathbf{P}_{i:}, \mathbf{C}_{i:} - \lambda_i \mathbf{1} - \boldsymbol{\mu} \rangle + \varepsilon \psi(\mathbf{P}_{i:}) \quad (15)$$

$$\stackrel{(*)}{=} \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \varepsilon > 0} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle + \varepsilon \left( \sum_i \psi^*((\mathbf{C}_{i:} - \lambda_i \mathbf{1} - \boldsymbol{\mu})/\varepsilon) - \eta \right). \quad (\text{Dual-ROT})$$

In  $(*)$  we have used that  $\psi^*(\mathbf{x}) = \sup_{\mathbf{y} \geq 0} \langle \mathbf{x}, \mathbf{y} \rangle - \psi(\mathbf{y})$ . From Danskin's theorem [10], one can recover the solution of the primal

$$\forall i, \mathbf{P}_{i:}^* = \nabla\psi^*((\mathbf{C}_{i:} - \lambda_i^* \mathbf{1} - \boldsymbol{\mu}^*)/\varepsilon^*). \quad (16)$$

Using matrix notations yields  $\mathbf{P}^* = \nabla\psi^*((\mathbf{C} - \boldsymbol{\lambda}^* \oplus \boldsymbol{\mu}^*)/\varepsilon^*)$  where  $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*, \varepsilon^*)$  are the solution of the dual problem (Dual-ROT).

## B.2 OT with Pointwise Constraints on Either Sources or Targets (OTARI-s and OTARI-t) : proof of Proposition 2

**Proposition 2.** Let  $(\mathbf{a}, \mathbf{b}, \xi)$  be such that  $\Pi(\mathbf{a}, \mathbf{b}) \cap \mathcal{B}_\psi(\xi)$  has an interior point and let  $\mathbf{P}^*$  solve

$$\min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle \quad \text{s.t.} \quad \mathbf{P} \in \mathcal{B}_\psi(\xi). \quad (\text{OTARI-s})$$

Let  $\varepsilon^*$  be the optimal dual variable associated with the constraint  $\mathbf{P} \in \mathcal{B}_\psi(\xi)$ . If  $\varepsilon^* > 0$ , then it holds  $\mathbf{P}^* = \text{Proj}_{\Pi(\mathbf{a}, \mathbf{b}) \cap \mathcal{B}_\psi(\xi)}^{D_\psi}(\mathbf{K}_{\varepsilon^*})$  for any  $0 < \varepsilon \leq \min_i \varepsilon_i^*$ . Moreover it holds



$$\mathbf{P}^* = \nabla \psi^* (\text{diag}(\boldsymbol{\varepsilon}^*)^{-1}(\mathbf{C} - \boldsymbol{\lambda}^* \oplus \boldsymbol{\mu}^*)) \text{ where } (\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*, \boldsymbol{\varepsilon}^*) \text{ solve the following dual}$$

$$\max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\varepsilon} > \mathbf{0}} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle + \langle \boldsymbol{\varepsilon}, \psi^* (\text{diag}(\boldsymbol{\varepsilon})^{-1}(\mathbf{C} - \boldsymbol{\lambda} \oplus \boldsymbol{\mu})) - \psi(\mathbf{e}_\xi) \mathbf{1} \rangle. \quad (\text{Dual-OTARI-s})$$

*Proof.* Again we break down the proof, focusing on the primal and then on the dual approach.

### Part I : Proof of the Bregman projection.

The proof is almost identical to the one in Proposition 1. We use the same notations for simplicity. The only difference brought by the pointwise constraint is that  $\boldsymbol{\nu}^*$  is now vectorial. The first-order optimality condition for the Bregman projection problem reads

$$\mathbf{C} + \sigma(\text{diag}(\boldsymbol{\nu}^*) + \mathbf{I}_{N_S})\nabla\psi(\mathbf{P}^*) - \sigma\boldsymbol{\lambda}^* \oplus \boldsymbol{\mu}^* - \sigma\boldsymbol{\Omega}^* = \mathbf{0}. \quad (17)$$

Again using the same notations as before, the first order KKT condition for problem (ROT) reads

$$\mathbf{C} + \text{diag}(\boldsymbol{\varepsilon}^*)\nabla\psi(\mathbf{T}^*) - \boldsymbol{\rho}^* \oplus \boldsymbol{\kappa}^* - \boldsymbol{\Lambda}^* = \mathbf{0}. \quad (18)$$

We end the proof by following the same reasoning as for Proposition 1, choosing for any  $i$ ,  $\nu_i = \varepsilon_i^*/\sigma - 1 \geq 0$ .

### Part II : Dual Problem of (OTARI-s).

The optimization problem (OTARI-s) writes

$$\min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle \quad \text{s.t.} \quad \mathbf{P} \in \mathcal{B}_\psi(\xi). \quad (19)$$

Introducing the dual variables  $\boldsymbol{\lambda} \in \mathbb{R}^n$  and  $\boldsymbol{\mu} \in \mathbb{R}^n$  for the marginals and  $\boldsymbol{\varepsilon} \in \mathbb{R}_+^n$  for the constraint  $\mathbf{P} \in \mathcal{B}_\psi(\xi)$ . The problem can be formulated as

$$\min_{\mathbf{P} \geq \mathbf{0}} \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\varepsilon} \geq \mathbf{0}} \langle \mathbf{P}, \mathbf{C} \rangle + \langle \boldsymbol{\lambda}, \mathbf{a} - \mathbf{P}\mathbf{1} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} - \mathbf{P}^\top \mathbf{1} \rangle + \langle \boldsymbol{\varepsilon}, \psi(\mathbf{P}) - \psi(\mathbf{e}_\xi) \mathbf{1} \rangle. \quad (20)$$

When  $\boldsymbol{\varepsilon}^* > \mathbf{0}$ , relying on strong duality to invert the min and max operators, the problem reduces to

$$\max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\varepsilon} > \mathbf{0}} \min_{\mathbf{P} \geq \mathbf{0}} \langle \mathbf{P}, \mathbf{C} \rangle + \langle \boldsymbol{\lambda}, \mathbf{a} - \mathbf{P}\mathbf{1} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} - \mathbf{P}^\top \mathbf{1} \rangle + \langle \boldsymbol{\varepsilon}, \psi(\mathbf{P}) - \psi(\mathbf{e}_\xi) \mathbf{1} \rangle \quad (21)$$

$$= \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\varepsilon} > \mathbf{0}} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle - \langle \boldsymbol{\varepsilon}, \psi(\mathbf{e}_\xi) \mathbf{1} \rangle + \min_{\mathbf{P} \geq \mathbf{0}} \langle \mathbf{P}, \mathbf{C} - \boldsymbol{\lambda} \mathbf{1}^\top - \mathbf{1} \boldsymbol{\mu}^\top \rangle + \langle \boldsymbol{\varepsilon}, \psi(\mathbf{P}) \rangle \quad (22)$$

$$= \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\varepsilon} > \mathbf{0}} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle - \langle \boldsymbol{\varepsilon}, \psi(\mathbf{e}_\xi) \mathbf{1} \rangle + \min_{\mathbf{P} \geq \mathbf{0}} \sum_i \langle \mathbf{P}_{i:}, \mathbf{C}_{i:} - \lambda_i \mathbf{1} - \boldsymbol{\mu} \rangle + \varepsilon_i \psi(\mathbf{P}_{i:}) \quad (23)$$

$$= \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\varepsilon} > \mathbf{0}} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle - \langle \boldsymbol{\varepsilon}, \psi(\mathbf{e}_\xi) \mathbf{1} \rangle + \sum_i \varepsilon_i \psi^*((\mathbf{C}_{i:} - \lambda_i \mathbf{1} - \boldsymbol{\mu})/\varepsilon_i) \quad (24)$$

$$\stackrel{(*)}{=} \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\varepsilon} > \mathbf{0}} \langle \boldsymbol{\lambda}, \mathbf{a} \rangle + \langle \boldsymbol{\mu}, \mathbf{b} \rangle + \langle \boldsymbol{\varepsilon}, \psi^*((\mathbf{C} - \boldsymbol{\lambda} \oplus \boldsymbol{\mu}) \circ \boldsymbol{\varepsilon} \mathbf{1}^\top) - \psi(\mathbf{e}_\xi) \mathbf{1} \rangle. \quad (\text{Dual-OTARI-s})$$

In  $(*)$  we used that  $\psi^*(\mathbf{X}) = (\psi^*(\mathbf{X}_{1:}), \dots, \psi^*(\mathbf{X}_{N:}))^\top$ . From Danskin's theorem [10], one can recover the solution of the primal

$$\mathbf{P}^* = \nabla \psi^* (\text{diag}(\boldsymbol{\varepsilon}^*)^{-1}(\mathbf{C} - \boldsymbol{\lambda}^* \oplus \boldsymbol{\mu}^*)) \quad (25)$$

where  $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*, \boldsymbol{\varepsilon}^*)$  solve (Dual-OTARI-s).  $\square$

## C $\psi$ -Bregman Projections for $\psi_{\text{KL}}$ and $\psi_2$

In this section, we detail the expressions of the projections used in the alternating Bregman projection approach.

### C.1 KL Projections

**Proposition 3.** When  $D_\psi$  is the KL divergence  $D_{\text{KL}}$ , one has for a matrix  $\mathbf{K} \in \mathbb{R}_+^{N_S \times N_T}$ ,

$$\text{Proj}_{\Pi(\mathbf{a}) \cap \mathcal{B}_{\text{KL}}(\xi)}^{\text{KL}}(\mathbf{K}) = \text{diag}(\boldsymbol{\Lambda} \mathbf{1})^{-1} \boldsymbol{\Lambda} \quad \text{with} \quad \boldsymbol{\Lambda} = \exp(\text{diag}(\mathbf{1} + \boldsymbol{\gamma}^*)^{-1} \log \mathbf{K}) \quad (26)$$

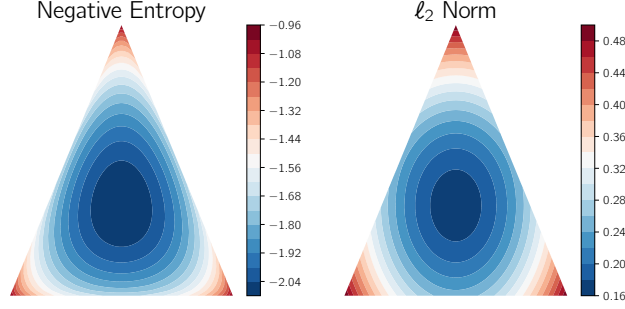


Figure 3:  $\sum_i \psi(p_i)$  plotted over the 3 dimensional probability simplex for  $\psi_{\text{KL}}$  (negative Shannon entropy) and  $\psi_2 : \mathbf{x} \rightarrow \frac{1}{2} \|\mathbf{x}\|_2^2$ . Unlike  $\psi_{\text{KL}}$ , the level sets of  $\psi_2$  intercept with the boundaries of the simplex thus leading to potentially sparse solutions when used to regularize OT.

where  $\gamma^* \geq \mathbf{0}$  is the optimal dual variable associated with the constraint  $\mathcal{B}_{\text{KL}}(\xi)$ .

*Proof.* The KL projection of a matrix  $\mathbf{K} \in \mathbb{R}_+^{N_S \times N_T}$  onto  $\Pi(\mathbf{a}) \cap \mathcal{B}_1(\xi)$  is the following problem.

$$\min_{\mathbf{P} \in \mathbb{R}_+^{N_S \times N_T}} \text{KL}(\mathbf{P}|\mathbf{K}) = \langle \mathbf{P}, \log(\mathbf{P} \oslash \mathbf{K}) - \mathbf{1}\mathbf{1}^\top \rangle \quad (27)$$

$$\text{s.t. } \forall i \in \llbracket N_S \rrbracket, \text{H}(\mathbf{P}_{i\cdot}) \geq \eta \quad (28)$$

$$\mathbf{P}\mathbf{1} = \mathbf{a}. \quad (29)$$

where for  $\mathbf{p} \in \mathbb{R}_+^{N_S}$ ,  $\text{H}(\mathbf{p}) = -\langle \mathbf{p}, \log \mathbf{p} - \mathbf{1} \rangle$  is the Shannon entropy. The associated Lagrangian writes

$$\mathcal{L}(\mathbf{P}, \boldsymbol{\lambda}, \boldsymbol{\gamma}) = \langle \mathbf{P}, \log \mathbf{P} - \log \mathbf{K} - \mathbf{1}\mathbf{1}^\top \rangle + \langle \boldsymbol{\gamma}, \boldsymbol{\eta}\mathbf{1} - \text{H}(\mathbf{P}) \rangle + \langle \boldsymbol{\lambda}, \mathbf{a} - \mathbf{P}\mathbf{1} \rangle. \quad (30)$$

Strong duality holds hence any optimal primal-dual variables  $(\mathbf{P}^*, \boldsymbol{\gamma}^*, \boldsymbol{\lambda}^*)$  must satisfy the KKT conditions. The first-order optimality condition gives

$$\nabla_{\mathbf{P}} \mathcal{L}(\mathbf{P}^*, \boldsymbol{\gamma}^*, \boldsymbol{\lambda}^*) = \log(\mathbf{P}^* \oslash \mathbf{K}) + \text{diag}(\boldsymbol{\gamma}^*) \log \mathbf{P}^* - \boldsymbol{\lambda}^* \mathbf{1}^\top = \mathbf{0}. \quad (31)$$

Isolating  $\mathbf{P}^*$  yields

$$\forall (i, j) \in \llbracket N_S \rrbracket \times \llbracket N_T \rrbracket, \quad P_{ij}^* = \frac{1}{u_i} \exp((\log K_{ij}) / (1 + \gamma_i^*)) \quad (32)$$

where  $u_i = \exp(-\lambda_i / (1 + \gamma_i^*))$ . Given the marginal constraint, we have

$$u_i = a_i^{-1} \sum_{j \in \llbracket N_T \rrbracket} \exp((\log K_{ij}) / (1 + \gamma_i^*)). \quad (33)$$

We are now left with  $\mathbf{P}^*$  as a function of  $\boldsymbol{\gamma}$ . Plugging  $\mathbf{P}^*$  in  $\mathcal{L}$  yields the dual function  $\boldsymbol{\gamma} \mapsto \mathcal{G}(\boldsymbol{\gamma})$ . This function is concave (property of the dual problem) and its gradient reads:

$$\nabla_{\boldsymbol{\gamma}} \mathcal{G}(\boldsymbol{\gamma}) = (\log \xi + 1) \mathbf{1} - \text{H}(\mathbf{P}^*(\boldsymbol{\gamma})). \quad (34)$$

Similarly to [25], one can show that the above gradient is canceled for a unique  $\bar{\boldsymbol{\gamma}}$ . The optimal dual variable is then given by  $\boldsymbol{\gamma}^* = [\bar{\boldsymbol{\gamma}}]_+$ .

□

## C.2 Euclidean Projections

For the Euclidean case, we break down the projection into  $\ell_2$  norm and marginal projections. Starting with the marginal projection, we have the following expression [18]

$$\text{Proj}_{\Pi(\mathbf{a})}^{\ell_2}(\mathbf{K}) = [\boldsymbol{\lambda}^* \mathbf{1}^\top + \mathbf{K}]_+ \quad (35)$$

where  $\boldsymbol{\lambda}^*$  is such that  $[\boldsymbol{\lambda}^* \mathbf{1}^\top + \mathbf{K}]_+ \in \Pi(\mathbf{a})$ .

Focusing on the  $\ell_2$  norm we have the following result.

**Proposition 4.** One has

$$\text{Proj}_{\mathcal{B}_2(\xi)}^{\ell_2}(\mathbf{K}) = \text{diag}(\boldsymbol{\gamma}^*)^{-1}\mathbf{K} \quad (36)$$

where for any  $i$ ,  $\gamma_i^* = \max(\xi^{1/2}\|\mathbf{K}_{i:}\|_2, 1)$ .

*Proof.* The  $D_2$  projection of a matrix  $\mathbf{K} \in \mathbb{R}_+^{N_S \times N_T}$  onto  $\mathcal{B}_2(\xi)$  reduces to

$$\min_{\mathbf{P} \in \mathbb{R}_+^{N_S \times N_T}} D_2(\mathbf{P}|\mathbf{K}_\varepsilon) = \frac{1}{2} \langle \mathbf{P}^{\odot 2}, \mathbf{1} \rangle - \langle \mathbf{P}, \mathbf{K} \rangle \quad (37)$$

$$\text{s.t. } \forall i \in \llbracket N_S \rrbracket, \|\mathbf{P}_{i:}\|_2^2 \leq (1/\xi). \quad (38)$$

Introducing the dual variable  $\boldsymbol{\omega} \in \mathbb{R}_+^n$ , the Lagrangian writes:

$$\mathcal{L}(\mathbf{P}, \boldsymbol{\omega}, \boldsymbol{\Omega}) = \frac{1}{2} \langle \mathbf{P}^{\odot 2}, \mathbf{1} \rangle - \langle \mathbf{P}, \mathbf{K} \rangle + \frac{1}{2} \sum_i \omega_i (\|\mathbf{P}_{i:}\|_2^2 - (1/\xi)). \quad (39)$$

$\mathbf{P}^*$  solves the primal problem if and only if there exists  $\boldsymbol{\omega}^*$  that satisfies the KKT conditions. The first-order condition yields

$$\nabla_{\mathbf{P}} \mathcal{L}(\mathbf{P}^*, \boldsymbol{\omega}^*, \boldsymbol{\Omega}^*) = -\mathbf{K} + \text{diag}(\boldsymbol{\omega}^* + \mathbf{1})\mathbf{P}^* = \mathbf{0}. \quad (40)$$

Hence it follows

$$\mathbf{P}^* = \text{diag}(\boldsymbol{\omega}^* + \mathbf{1})^{-1}\mathbf{K}. \quad (41)$$

To satisfy the KKT conditions, one has to find the root  $\boldsymbol{\omega}^*$  of the following independent problems

$$\forall i, (\omega_i + 1)^2 = \xi \|\mathbf{K}_{i:}\|_2^2. \quad (42)$$

Thus we have  $\omega_i + 1 = \xi^{1/2}\|\mathbf{K}_{i:}\|_2$  and taking  $\omega_i \geq 0$  into account yields the result.  $\square$