

---

# A PROXIMAL-SINKHORN-NEWTON METHOD FOR ENTROPIC OPTIMAL TRANSPORT

Anonymous authors

Paper under double-blind review

## ABSTRACT

Entropic optimal transport (OT) enables efficient distribution alignment through the Sinkhorn method. However, it suffers from numerical instability and slow convergence under weak entropic regularization. We propose a two-stage framework that establishes an inexact-to-exact paradigm to address these challenges. The first stage employs an inexact proximal point method to decompose the entropic OT into simpler subproblems, yielding approximate solutions with superior numerical stability. The second stage employs a sparse Newton method with global convergence and a locally quadratic rate to refine the approximate solutions. Compared to previous Newton-based algorithms, it accelerates updates and prevents the objective score-value from plateauing during optimization. With numerical instability handled in the first stage, Sinkhorn iterations-scaling can provide an alternative to the Newton method under relatively heavy entropic regularization. The yielding-resulting Proximal-Sinkhorn-Newton method enjoys the strengths of three approaches and outperforms the baselines across various regularizations and error tolerances, achieving over 20× speedups on some problems compared to other Newton-based methods.

## 1 INTRODUCTION

Optimal transport (OT) calculates-seeks the best transportation plan from an ensemble of sources to targets (Linial et al., 1998; Peyré et al., 2017). It is an important task in machine learning, such as GAN (Arjovsky et al., 2017; Genevay et al., 2018), graph matching (Shen et al., 2024), domain adaption-adaptation (Nguyen et al., 2021; Turrisi et al., 2022), statistical learning (Oneto et al., 2020; Huynh et al., 2020), robust and efficient unsupervised learning (Lei et al., 2019; Onken et al., 2021), feature selection (Cao and Zhang, 2024; Cao et al., 2026), and so on. For an overview of further applications of optimal transport, readers are referred to (Peyré and Cuturi, 2019).

In this work, we focus on the entropic optimal transport problem of the form

$$P_\beta = \arg \min_{P \in \mathbb{U}} C \cdot P + \frac{1}{\beta} \mathcal{H}(P), \quad \mathcal{H}(P) := P \cdot \mathring{\log} P \quad (1)$$

where  $P \in \mathbb{R}_+^{n \times n}$  represents a transport plan,  $\mathbb{U} := \{P \in \mathbb{R}_+^{n \times n} : P\mathbf{1} = \mathbf{a}, P^\top \mathbf{1} = \mathbf{b}\}$  is the feasible region where  $\mathbf{a}, \mathbf{b} \in \mathbb{R}_+^n$  are the source and target distributions respectively,  $C \in \mathbb{R}^{n \times n}$  is a cost matrix,  $C \cdot P = \sum C_{ij} P_{ij}$  is the transport cost,  $\beta > 0$  is the regularization parameter,  $\mathcal{H}(\cdot)$  is the negative entropy, and  $\mathring{\log}$  denotes the element-wise logarithm operator i.e.,  $P \cdot \mathring{\log} P = \sum P_{ij} \log(P_{ij})$ . Cuturi (2013) demonstrates that the solution to (1) admits a closed-form characterization:

$$P_\beta = \Gamma_{\mathbb{U}}^{\text{kl}} \left( \mathring{\exp}(-\beta C) \right) = D_{(\mathbf{u})} \mathring{\exp}(-\beta C) D_{(\mathbf{v})}, \quad (2)$$

where  $\Gamma_{\mathbb{U}}^{\text{kl}}(\cdot)$  denotes the projection onto the set  $\mathbb{U}$  with respect to the Kullback–Leibler (KL) divergence (Peyré et al., 2017),  $\mathring{\exp}(\cdot)$  represents element-wise exponential operator,  $D_{(\mathbf{u})}$  and  $D_{(\mathbf{v})}$  are diagonal matrices constructed from positive scaling vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}_+^n$  which. The two vectors can be computed iteratively via the classical Sinkhorn algorithm (Sinkhorn, 1964; Cuturi, 2013):

$$K = \mathring{\exp}(-\beta C), \quad \mathbf{v}^{(l+1)} = \mathbf{a} \oslash (K \mathbf{u}^{(l)}), \quad \mathbf{u}^{(l+1)} = \mathbf{b} \oslash (K^\top \mathbf{v}^{(l+1)}). \quad (3)$$

where  $\oslash$  is the element-wise division.

According to (2), solving the entropic OT relies on two components: the kernel matrix  $K$  and the scaling vectors  $\mathbf{u}, \mathbf{v}$ . These introduce two intrinsic computational challenges for entropic OT with a large  $\beta$ : (a) numerical instability caused by division by zero, which can occur when  $K$  contains an entire row of **numeric**-zeros; and (b) the iterative computation of  $\mathbf{u}$  and  $\mathbf{v}$  with complexity tied to **the**  $\beta$ . These challenges are inherent to the entropic OT and continue to motivate ongoing research.

**Numerical instability.** The stable Sinkhorn (Benamou et al., 2015) addresses the numerical instability by performing computations in log-space, but introduces considerable overhead. The proximal point method (POT) (Xie et al., 2020) generates a sequence  $\{P^{(t)}\}_{t=1,2,\dots}$  to obtain the solution **of**-to OT by solving a series of entropic OT problems with small values of  $\beta$ . To improve efficiency, the inexact proximal point method (IPOT) (Xie et al., 2020) computes the proximal operator approximately at each iteration, yet  $\{\hat{P}^{(t)}\}_{t=1,2,\dots}$  still converges to the exact solution (shown in Figure 1). However, the algorithm exhibits only linear convergence in OT and does not benefit from the efficiency of entropic OT, thus becoming potentially expensive for large-scale problems.

**Computational complexity.** The Sinkhorn algorithm enjoys an  $\mathcal{O}(n^2)$  per-iteration cost, but due to its sub-linear convergence, **attaining a moderately high accuracy may still demand the algorithm may require** a prohibitively large number of iterations to attain a moderately high accuracy (Tang and Qiu, 2024). The Newton method has a local quadratic convergence rate, while even a single Newton step requires  $\mathcal{O}(n^3)$  cost (Brauer et al., 2017). The Sinkhorn-Newton-Sparse (SNS) algorithm (Tang et al., 2024) sparsifies the Hessian matrix to reduce the computational cost to  $\mathcal{O}(n^2)$ . However, its warm start may require hundreds of Sinkhorn iterations (as shown in Figure 1). Tang and Qiu (2024) propose **a**-an SSNS to overcome this problem, whereas the objective **score**-value may experience a period of stagnation during the Newton iterateiteration.

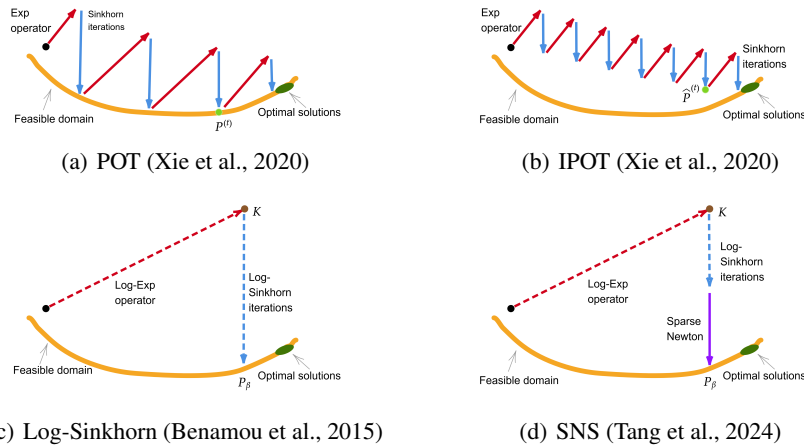


Figure 1: Schematic of the path of methods. The distance is in the Bregman sense. POT and IPOT are designed for OT. In fact, they require far more iterations than those illustrated in the figure to converge; log-Sinkhorn and SNS are designed for entropic OT.

In light of the above challenges, a natural question arises:

How to design an entropic optimal transport algorithm that maintains stability and efficiency under varying entropic regularization parameters and error tolerances?

A promising approach is to weave complementary algorithms into a single framework that leverages their respective strengths. Specifically, we employ the inexact proximal point method to stably construct a transport plan. If this plan does not meet the desired accuracy, it can **serve**-serve as a high-quality initializer for faster secondary solvers. Such an idea is supported by the following theorem.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

**Theorem 1** *Inexact-to-Exact Transition*

*Let  $\hat{P}^{(t)}$  denote the transport plan updated at the  $t$ -th iteration in the inexact proximal point method with a fixed proximal parameter  $\Delta\beta$ , then*

$$P_\beta = \Gamma_{\mathbb{U}}^{\text{kl}}(\hat{P}^{(t)}), \text{ where } \beta = t\Delta\beta. \tag{4}$$

This theorem establishes a fundamental connection between the inexact proximal point method and the entropic OT. Specifically, it allows the inexact iterate  $\hat{P}^{(t)}$  to be refined through the KL-projection into the exact solution of the entropic OT problem. Interpreting  $\hat{P}^{(t)}$  as a processed kernel matrix reveals two distinct benefits over the naive counterpart  $K$ : improved numerical robustness and reduced distance to the solution  $P_\beta$  (shown in Figure 1). These properties make it a more desirable input for subsequent refinement algorithms such as the Sinkhorn [iterations-scaling](#) or the Newton method.

The contributions of this paper are threefold:

- **Numerical instability: addressed via an inexact proximal point framework.** We use a proximal point method that safely decomposes the problem into a sequence of easier and safer entropic OT problems. By employing an inexact proximal operator, we derive an Inexact Proximal point [Method-method](#) for Entropic OT (IP-EOT), stably and efficiently producing an approximate solution that other methods can further refine.
- **Computational complexity: reduced via a sparse Newton method.** We introduce a novel sparse Newton method with global convergence and locally quadratic rates. The positive definiteness of the sparsified Hessian matrix [enabling-enables](#) the use of efficient incomplete Cholesky conjugate gradient solvers to accelerate each Newton update. Furthermore, we integrate line search with Sinkhorn [iterations-scaling](#) to prevent the objective [score-value](#) from plateauing during the optimization process.
- **Adaptive hybrid solver.** Our framework leverages IP-EOT either as a low-accuracy standalone solver or as a high-accuracy warm-start for both Sinkhorn and Newton methods, thereby overcoming instability of the Sinkhorn algorithm and expensive initialization of the Newton method. A proposed adaptive criterion selects the secondary solver by exploiting the Hessian’s sparsity dependence on  $\beta$ : the Newton method is preferred for large  $\beta$  (sparser Hessian), while the Sinkhorn [iterations-scaling](#) is chosen for smaller  $\beta$ .

In short, we address two fundamental challenges in entropic optimal transport—numerical instability and high computational complexity—by introducing a unified inexact-to-exact framework, termed Proximal-Sinkhorn-Newton (PSN), which intelligently integrates complementary algorithms. The overall architecture of PSN is illustrated in Figure 2.

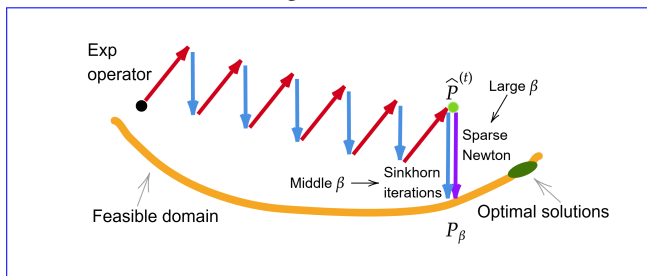


Figure 2: The trajectory of PSN consists of two stages: (a) IP-EOT for [an](#) efficient and stable approximation  $\hat{P}^{(t)}$ , and (b) KL-projection by either Sinkhorn [iterations-scaling](#) or the Newton method for refinement.

## 2 RELATED WORK

**Computational Optimal Transport.** Since the Sinkhorn algorithm (Cuturi, 2013) made entropic optimal transport practical, researchers have polished it from every angle. First-order alternating

minimisation (Guminov et al., 2021), Greenkhorn (Altschuler et al., 2017) and their greedy variants — speed up the basic matrix-scaling loops. Adaptive primal-dual strategies, such as the similar-triangle method (Dvurechensky et al., 2020) and accelerated mirror descent (Li et al., 2025), dynamically tune step-sizes to reduce the iteration count. When high accuracy is required, Newton-type solvers exploit second-order information to yield local quadratic convergence (Brauer et al., 2017; Tang et al., 2024; Tang and Qiu, 2024). Yet, an efficient method applicable to a wide range of  $\beta$  is still lacking, which motivates this work.

**Newton method.** The Newton method achieves rapid convergence, whereas handling a dense Hessian incurs  $\mathcal{O}(n^3)$  complexity (Powell and Toint, 1979). To mitigate this, [Tang et al. \(2024\); Tang and Qiu \(2024\); Tang et al. \(2024\); Tang and Qiu \(2024\); Kemertas et al. \(2025\)](#) sparsify the Hessian matrix to reduce the per-iteration cost to approximately  $\mathcal{O}(n^2)$ . Iterative solvers such as Conjugate Gradient (CG) (Golub and Van Loan, 2013) benefit significantly from preconditioning (Caraba, 2008), yet this has been overlooked in prior work on optimal transport. We therefore apply a preconditioner to accelerate the Newton method.

**Proximal point method.** The proximal point method transforms a complex optimization problems into a sequence of simpler subproblems by introducing a proximal operator (Rockafellar, 1976). Based on the Bregman distances, Chen and Teboulle (1993) formalized the Bregman proximal point algorithm and established its convergence. Recently, this method ~~have~~ has been applied to OT, substantially reducing computational complexity (Xie et al., 2020; Chu et al., 2023; Chizat, 2024). The unclear trajectory of the regularization parameter during iteration has hindered use of the proximal point method in entropic OT, which our work seeks to address.

### 3 PROXIMAL POINT METHODS FOR ENTROPIC OPTIMAL TRANSPORT

We propose an Inexact Proximal point method for ~~the~~ Entropic OT (IP-EOT) that yields an approximation  $\hat{P}_\beta$  ~~(interpretable as result from insufficient Sinkhorn iterations) by making by explicitly incorporating  $\beta$  explicit in into~~ the proximal step. ~~We further show This approximation can also be derived from an approximate KL-projection  $\hat{\Gamma}_U^{\text{kl}}(\exp(-\beta C))$ , a typical instance of which is early termination of Sinkhorn scaling. Furthermore, we demonstrate that  $\hat{P}_\beta$  can be refined into the exact solution  $P_\beta$  via the through an exact KL-projection.~~

**Proximal point methods** are iterative algorithms for convex optimization, well-suited for nonsmooth or composite objectives. Their core principle involves solving a regularized subproblem at each step. This process optimizes the objective while ensuring the stability of the iterative sequence. To solve the OT (without the entropic term), the proximal point method reads as [follows](#)

$$P^{(t+1)} = \arg \min_{P \in \mathbb{U}} C \odot P + \beta^{(t+1)} \mathcal{B}(P, P^{(t)}), \quad (5)$$

$$\mathcal{B}(P, P^{(t)}) = P \cdot \log(P \oslash P^{(t)}) - (P - P^{(t)}) \cdot \mathbf{1}\mathbf{1}^\top, \quad (6)$$

~~where  $\odot$  is the element-wise product.~~  $\mathcal{B}(\cdot, \cdot)$  is the Bregman divergence used to measure the distance between two probability vectors/matrices (Benamou et al., 2015). This formula can be equivalently rewritten as  $\div$

$$P^{(t+1)} = \arg \min_{P \in \mathbb{U}} C^{(t)} \odot P + \beta^{(t+1)} \mathcal{H}(P), \quad C^{(t)} = C - \beta^{(t+1)} \log P^{(t)}. \quad (7)$$

The update can be computed by applying a KL-projection to a modified kernel matrix, yielding:

$$P^{(t+1)} = \Gamma_U^{\text{kl}} \left( \exp(-C\beta^{(t+1)}) \odot P^{(t)} \right). \quad (8)$$

As  $t \rightarrow \infty$ ,  $P^{(t)}$  converges to the optimal transport plan. To improve efficiency, Xie et al. (2020) ~~propose~~ proposed the inexact proximal point method (IPOT), which performs only one Sinkhorn iteration to approximate the projection in (8). The resulting sequence  $\{\hat{P}^{(t)}\}_{t=1,2,\dots}$  can efficiently converge to the optimal transport plan.

**Inexact proximal point method for entropic OT.** Theorem 1 allows the inexact proximal point method to efficiently offer a  ~~$\hat{P}_\beta$  for the entropic optimal transport to compute  $P_\beta$ . The  $\hat{P}_\beta$ .~~ This

theorem relies on explicitly tracking the regularization parameter  $\beta$  across proximal iterations, as detailed below.

**Lemma 1 (Explicit form of  $\beta$ )** *The iterative schemes of the exact/inexact proximal point method for entropic OT admit the following explicit representation:*

$$P^{(t)} = \Gamma_{\mathbb{U}}^{\text{kl}}\left(\bigodot_{i=1}^t \text{exp}(-C\beta^{(i)})\right), \quad \hat{P}^{(t)} = \hat{\Gamma}_{\mathbb{U}}^{\text{kl}}\left(\bigodot_{i=1}^t \text{exp}(-C\beta^{(i)})\right). \quad (9)$$

$\bigodot$  represents element-wise multiplication of a series of matrices. In particular, when each  $\beta^{(i)}$  is fixed as  $\Delta\beta$ , the regularization  $\beta$  satisfies  $\beta = t\Delta\beta$ , which yields the following explicit expression:

$$\hat{P}^{(t)} = \hat{\Gamma}_{\mathbb{U}}^{\text{kl}}(\text{exp}(-t\Delta\beta C)) = \hat{P}_{\beta}, \quad \beta = t\Delta\beta. \quad (10)$$

~~where  $\hat{\Gamma}_{\mathbb{U}}^{\text{kl}}$  represents an approximate KL-projection~~ The  $\hat{P}_{\beta}$  can admit the form  $\hat{P}_{\beta} = D_{(\frac{1}{\beta}\mathbf{1})} \text{exp}(-\beta C) D_{(\frac{1}{\beta}\mathbf{1})}$  and satisfies

$$\rho_{\mathbb{U}}(P_{\beta}^0) > \rho_{\mathbb{U}}(\hat{P}_{\beta}) > 0, \quad (11)$$

where  $\rho_{\mathbb{U}}(P) = \|P\mathbf{1} - \mathbf{a}\|_1 + \|P^{\top}\mathbf{1} - \mathbf{b}\|_1$  represents the marginal error and  $P_{\beta}^0 = D_{(\frac{1}{\beta}\mathbf{1})} \text{exp}(-\beta C) D_{(\frac{1}{\beta}\mathbf{1})}$  is a classical initialization for KL-projections. The next Lemma establishes a theoretical bridge between the inexact transport plan and the exact transport plan.

**Lemma 2 (Inexact to Exact)** *Let  $\hat{P}_{\beta}$  be an approximation to ~~the~~  $P_{\beta}$ , ~~then~~*

$$P_{\beta} = \Gamma_{\mathbb{U}}^{\text{kl}}(\hat{P}_{\beta}). \quad (12)$$

This result enables the inexact approximation  $\hat{P}_{\beta}$  to serve as a warm start for computing the exact solution  $P_{\beta}$  via the KL-projection.

Building on the above results, we propose an IP-EOT to provide an approximation  $\hat{P}_{\beta}$ , as outlined in Algorithm 1. The algorithm avoids the risk of numerical underflow by dividing the total regularization  $\beta$  into  $\ell$  steps (line 2) and progressively constructing the kernel matrix (line 6). ~~This strategy greatly enhances numerical stability, which can be further refined by adjusting  $\Delta\beta$ .~~ Our matrices  $\Delta K$  and  $\hat{P}$  are much less prone to numerical underflow than the Sinkhorn kernel matrix  $\text{exp}(-\beta C)$ . This advantage becomes more pronounced with a smaller  $\beta$ , leading to substantially improved numerical stability. With the aid of the IC preconditioner (line 4), line 5 employs a preconditioned conjugate gradient (PCG) method for computing  $\Delta z$ . Since the cost of each iteration is  $\mathcal{O}(n^2)$ , the computational complexity of the algorithm is  $\mathcal{O}(\ell n^2)$ .

---

#### Algorithm 1 IP-EOT

---

**Require:** ~~Probabilities~~ Probability vectors  $\{\mathbf{a}, \mathbf{b}\}$  with length  $m$  and  $n$ , cost matrix  $C$ , entropic parameter  $\beta$ , proximal parameter  $\Delta\hat{\beta}$

- 1:  $\ell = \lceil \frac{\beta}{\Delta\hat{\beta}} \rceil$
- 2:  $\Delta\beta \leftarrow \beta/\ell$
- 3:  $\Delta K \leftarrow \text{exp}(-\Delta\beta C)$
- 4:  $\hat{P} \leftarrow \mathbf{1}_m \mathbf{1}_n^{\top} / (mn)$
- 5: **for**  $t = 1, 2, \dots, \ell$
- 6:      $K \leftarrow \Delta K \odot \hat{P}$
- 7:      $\mathbf{u} \leftarrow \mathbf{a} \odot (K\mathbf{v}), \mathbf{v} \leftarrow \mathbf{b} \odot (K^{\top}\mathbf{u})$
- 8:      $\hat{P} \leftarrow \text{diag}(\mathbf{u})K\text{diag}(\mathbf{v})$
- 9: **end**
- 10: **return**  $\hat{P}$

---



---

#### Algorithm 2 Fast Sparse Newton

---

**Require:** Initialized transport plan  $P_0$ , Probabilities Probability vectors  $\{\mathbf{a}, \mathbf{b}\}$ ,  $C$ , marginal error threshold  $\rho_{th}, \beta$

- 1: **while**  $\rho > \rho_{th}$  **do**
- 2:     Compute  $H, \nabla f, \rho$  with (15)
- 3:      $\hat{H} \leftarrow \text{Sparsify}(H, \rho)$
- 4:      $L \leftarrow \text{IC}(\hat{H})$       $\triangleright \# \hat{H} \approx LL^{\top}$
- 5:      $\Delta z \leftarrow \text{PCG}(\hat{H}, -\nabla f, L)$
- 6:     **if**  $\alpha$  is searched **then**  $z \leftarrow z + \alpha\Delta z$
- 7:     **else** Update  $z$  by Sinkhorn iterations scaling
- 8:     Compute marginal error  $\rho$
- 9: **end while**
- 10:  $(\mathbf{x}, \mathbf{y}) \leftarrow z$
- 11:  $P \leftarrow \text{exp}(\beta(-C + \mathbf{x}\mathbf{1}^{\top} + \mathbf{1}\mathbf{y}^{\top}) - 1)$
- 12: **return**  $P$

---

Since the iteration number of IP-EOT is tied to  $\Delta\beta$ , the algorithm cannot adaptively stop based on marginal error. Thus, it should be augmented with alternative solvers, switching to Sinkhorn or Newton methods if the IP-EOT solution yields insufficient accuracy.

#### 4 FAST SPARSE NEWTON METHOD FOR ENTROPIC OPTIMAL TRANSPORT

The Newton method is based on the Lagrange function  $\mathcal{L}$  associated with the entropic OT (1):

$$\mathcal{L}(P, \mathbf{x}, \mathbf{y}) := \frac{1}{\beta} P \cdot \log P + C \cdot P - \mathbf{x} \cdot (P\mathbf{1} - \mathbf{a}) - \mathbf{y} \cdot (P^\top \mathbf{1} - \mathbf{b}). \quad (13)$$

Using the primal-dual and strong duality theory (Nocedal and Wright, 1999), one only needs to minimize the function  $f$ :

$$f(\mathbf{z}) = f(\mathbf{x}, \mathbf{y}) = - \min_{P \in \mathbb{U}} \mathcal{L}(P, \mathbf{x}, \mathbf{y}),$$

which is obtained by substituting  $P$  with  $\exp(\beta(-C + \mathbf{x}\mathbf{1}^\top + \mathbf{1}\mathbf{y}^\top) - 1)$ , to get the solution of (1):

$$f(\mathbf{z}) = f(\mathbf{x}, \mathbf{y}) = - \min_{P \in \mathbb{U}} \mathcal{L}(P, \mathbf{x}, \mathbf{y}), \quad (14)$$

The explicit form of  $f(\mathbf{z})$  follows by substituting  $P = \exp(\beta(-C + \mathbf{x}\mathbf{1}^\top + \mathbf{1}\mathbf{y}^\top) - 1)$  into (13). The corresponding gradient and Hessian ~~matrix~~ ~~matrices~~ of  $f$  are

$$\begin{aligned} \nabla_{\mathbf{z}} f(\mathbf{z}) &= (\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), \nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})) = (P\mathbf{1} - \mathbf{a}, P^\top \mathbf{1} - \mathbf{b}), \\ H(\mathbf{z}) = \nabla^2 f(\mathbf{z}) = \nabla^2 f(\mathbf{x}, \mathbf{y}) &= \beta \begin{bmatrix} \text{diag}(P\mathbf{1}) & P \\ P^\top & \text{diag}(P^\top \mathbf{1}) \end{bmatrix}. \end{aligned} \quad (15)$$

If  $H(\mathbf{z}_k)$  is invertible, the Newton method updates  $\mathbf{z}_{k+1}$  via the formula

$$\mathbf{z}_{k+1} = \mathbf{z}_k + \Delta \mathbf{z}_k, \quad \Delta \mathbf{z}_k = H(\mathbf{z}_k)^{-1} (-\nabla_{\mathbf{z}} f(\mathbf{z}_k)). \quad (16)$$

To avoid costly matrix inversion,  $\Delta \mathbf{z}$  is obtained by solving a linear system  $H\Delta \mathbf{z} = -\nabla f$ , which is commonly solved by the Conjugate Gradient (CG) method (Tang et al., 2024). However, a key challenge arises from the  $\mathcal{O}(n^3)$  computational cost per Newton iteration and the non-invertibility of the Hessian matrix  $H$ .

##### 4.1 FAST POSITIVE-DEFINITE-PRESERVING SPARSE SCHEME

Prior research has made substantial progress in characterizing complexity and irreversibility by sparsifying the Hessian. Tang et al. (2024) observed that, after some Sinkhorn iterations, the Hessian matrix can be well-approximated by a sparse structure. Building on this observation, they reduced the computational complexity of the single Newton update (16) from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(n^2)$  by enforcing a predetermined sparsity pattern on the Hessian matrix. ~~Building on this~~ ~~Following this direction~~, Tang and Qiu (2024) proposed an improved sparsification strategy that ensures the invertibility of the sparsified Hessian matrix. However, both approaches suffer from a significant ~~increasing~~ ~~increase~~ ~~in~~ computation time on certain large-scale problems. To overcome this limitation, we propose an efficient sparsification scheme that ~~grantees~~ ~~guarantees~~ an invertible sparse approximation of the Hessian matrix  $H$ .

The ~~sparsified matrix~~  ~~$\hat{H}$~~  is constructed by enforcing sparsity on the anti-diagonal submatrices  $P$  and  $P^\top$ :

$$\hat{H} = \beta \begin{bmatrix} \text{diag}(P\mathbf{1}) & \hat{P} \\ \hat{P}^\top & \text{diag}(P^\top \mathbf{1}) \end{bmatrix}. \quad (17)$$

The sparsification procedure comprises the following ~~steps~~: **Adaptive thresholding**: We truncate the matrices  $P$  and  $P^\top$  in  $H$  by zeroing out all elements with values smaller than a threshold

$\tau = \|\nabla f\|_1 = \rho$ , two steps:

$$\hat{P}_{ij}^0 = \begin{cases} P_{ij}, & P_{ij} \geq \tau, \\ 0, & P_{ij} < \tau, \end{cases} \quad \hat{P}_{ij} = \begin{cases} 0, & j = \arg \min_{k: \hat{P}_{ik}^0 > 0} \hat{P}_{ik}^0, (P - \hat{P}^0)_i = \mathbf{0}, \\ 0, & i = \arg \min_{k: \hat{P}_{kj}^0 > 0} \hat{P}_{kj}^0, (P - \hat{P}^0)_j = \mathbf{0}, \\ \hat{P}_{ij}^0, & \text{otherwise.} \end{cases} \quad (18)$$

where  $\rho$  is the marginal error, while capping  $\tau$  at  $10^{-4} \tau = \|\nabla f\|_1$  is guided by the objective of achieving quadratic convergence (detailed in the proof of Theorem 2). The second step, diagonal dominating operation, ensures that all rows of  $\hat{P}$  and  $\hat{P}^\top$  contain truncated elements, which also implies that  $\hat{H}$  exhibits diagonal dominance and positive definiteness. Appendix H presents the relationship between the sparsity level and  $\beta$ , and describes how to avoid excessive sparsity (discussed in Appendix H) sparsification which may slow down the optimization.

**Diagonal dominance enforcement:** If every row and column of the sparsified  $P$  contains truncated elements, the resulting Excessive sparsification, which removes important entries from the Hessian, may slow down the optimization. To avoid this, we reduce  $\tau$  whenever the number of nonzeros in  $\hat{P}$  falls below  $2n$ , ensuring that  $\hat{P}$  retains at least  $2n$  entries. The choice of  $2n$  is motivated by the fact that the solution of the corresponding OT contains at most  $2n - 1$  nonzero entries. Thus, if  $\hat{P}$  contains more than  $2n$  nonzeros, the corresponding  $\hat{H}$  is positive definite and hence invertible. Otherwise, we subtract the smallest entry from those rows and columns, ensuring diagonal dominance of symmetric  $\hat{H}$ , thereby guaranteeing its positive definiteness. The concise design, together with its  $\mathcal{O}(n^2)$  complexity, contributes to the remarkable scalability; the positive definiteness of  $\hat{H}$  enables stable and efficient CG-based Newton updates; adaptive threshold contributes to local quadratic convergence of the algorithm can be regarded as reliable. The detailed procedure can be found in the Appendix H.

**Theorem 2 (Local Quadratic Convergence)** Let the sequence  $\{z_k\}$  be generated by

$$z_{k+1} = z_k - \hat{H}_k^{-1} \nabla f(z_k), \quad (19)$$

where  $\hat{H}_k$  is sparsified from the Hessian matrix  $H_k = H(z_k)$  by the fast positive-definite-preserving sparse scheme. The sequence  $\{z_k\}$  converges locally quadratically to the minimizer  $z^*$ .

## 4.2 PRECONDITIONER FOR NEWTON METHOD

To enhance the efficiency of each Newton iteration, we employ an incomplete Cholesky (IC) preconditioner to concentrate the eigenvalue distribution of the Hessian matrix, thereby reducing the number of CG iterations. This design leverages the fact that CG converges more rapidly when the eigenvalues of  $H$  are tightly clustered (Hestenes et al., 1952) (see Appendix for details). The effectiveness of IC is demonstrated in Figure 9, which compares it with the diagonal preconditioner used in previous Newton-based algorithms for entropic OT (Brauer et al., 2017; Tang and Qiu, 2024).

## 4.3 LINE SEARCH WITH SINKHORN FALLBACK

The step size is typically determined by a line search procedure, such as backtracking, to ensure that the objective value decreases at each Newton iteration. In a backtracking line search, the step size is initialized to  $\alpha = 1$  and then repeatedly reduced by multiplying it with a constant factor  $\in (0, 1)$  (e.g., 0.5) until  $f(z_k + \alpha \Delta z_k) < f(z_k)$ . However, a practical issue arises when the line search yields an extremely small step size, causing the algorithm to enter a plateau phase in which the score barely decreases despite a considerable number of iterations. Tang and Qiu (2024) mitigated this issue by progressively strengthening the diagonal of the Hessian matrix, which shortens but does not completely eliminate the stagnation period.

To further address this problem, we propose a novel strategy. If the step size  $\alpha$  obtained from a backtracking line search fails to produce an improvement in the objective value

$$f(z_k + \alpha \Delta z_k) > f(z_k), \quad \forall \alpha \in \{1, 0.5, 0.5^2, 0.5^3, 0.5^4\}, \quad (20)$$

then the algorithm employs the Sinkhorn scaling to refine the solution. This can generate a more robust direction  $\mathbf{p}_k$  to escape the stagnation region and resume effective progress

$$f(\mathbf{z}_k + \mathbf{p}_k) < f(\mathbf{z}_k). \quad (21)$$

Once a suitable step size  $\alpha$  is found, the method resumes Newton updates, achieving both robustness and rapid convergence.

The robustness of this first-order direction is underpinned by its satisfaction of the Armijo and Wolfe conditions, as formalized in the following lemma. This not only guarantees sufficient descent but also prevents continued stagnation, thereby re-establishing the conditions necessary for the subsequent Newton steps to resume their rapid convergence.

**Lemma 3** *Let  $f(\mathbf{z})$  be defined in equation (14), and the sequence of search directions  $\{\mathbf{p}_k\}$  be generated by Sinkhorn iterations. Then in each update, direction  $\mathbf{p}_k$  with unit step size  $\alpha_k = 1$  satisfies the following conditions:*

$$f(\mathbf{z}_k + \mathbf{p}_k) \leq f(\mathbf{z}_k) + c_1 \nabla f(\mathbf{z}_k)^\top \mathbf{p}_k, \quad (22)$$

$$\nabla f(\mathbf{z}_k + \mathbf{p}_k)^\top \mathbf{p}_k \geq c_2 \nabla f(\mathbf{z}_k)^\top \mathbf{p}_k, \quad (23)$$

for constants  $0 < c_1 < c_2 < 1$ .

#### 4.4 FAST SPARSE NEWTON METHOD

Algorithm 2 outlines the fast sparse Newton method for entropic OT. In each iteration, the Hessian is sparsified (line 3) and preconditioned using an incomplete Cholesky (IC) factorization  $L$  (line 4). The main computational cost lies in computing  $\Delta \mathbf{z}$ , which requires  $\mathcal{O}(n \|\hat{H}\|_0)$  operations, where  $\|\hat{H}\|_0$  denotes the number of nonzero elements (line 5). At line 6, the step size  $\alpha$  is selected via a backtracking line search over  $\{0.5, 0.5^2, 0.5^3, 0.5^4\}$ . If none of these values improves the objective, we apply the Sinkhorn [iterations-scaling](#) for refinement. This adjustment helps Newton’s method to escape flat regions—, where progress would otherwise stall due to vanishing step sizes—, and is critical for establishing global convergence (line 7).

**Theorem 3 (Global Convergence)** *The sequence  $\{\mathbf{z}_k\}$  generated by Algorithm 2 converges globally to the minimizer  $\mathbf{z}^*$ .*

### 5 PROXIMAL-SINKHORN-NEWTON [METHOD](#)

This section presents the main algorithm, the Proximal-Sinkhorn-Newton (PSN), in Algorithm 3. The algorithm begins with the IP-EOT to produce an approximate solution  $\hat{P}_\beta$ . If ~~the~~  $\hat{P}_\beta$  meets the predefined criteria, it is returned as the final output. Otherwise, the algorithm refines the solution by switching to either the Sinkhorn [iterations-scaling](#) or the fast sparse Newton method, depending on the sparsity of  $\hat{H}$ . If  $\hat{H}$  contains fewer than  $\lambda n$  non-zero elements, the fast sparse Newton method is employed; otherwise, the Sinkhorn [iterations-are-scaling is](#) used. This approach is motivated by the relationship between sparsity and computational efficiency: Newton’s method benefits from sparser matrices, which occur as  $\beta$  increases (Tang et al., 2024). Thus, the strategy automatically favors Newton steps for larger  $\beta$  and Sinkhorn [iterations-scaling](#) for smaller  $\beta$ .

**Theorem 4** *Algorithm 3 converges globally to the solution of the corresponding entropic OT.*

## 6 NUMERICAL EXPERIMENTS

We evaluate the proposed algorithm PSN and other contributions from the following three parts:

- Q1. How does the performance of IP-EOT compare to the Sinkhorn algorithm?
- Q2. What are the advantages of using an IC preconditioner in the Newton method?
- Q3. [Whether PSN is scalable with respect to problem size and regularization strength](#)How does PSN perform across different problem sizes and regularization strengths?

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

---

**Algorithm 3** Proximal-Sinkhorn-Newton (PSN)

---

**Require:** ~~Probabilities~~ Probability vectors  $\{\mathbf{a}, \mathbf{b}\}$ , cost matrix  $C$ , entropic parameter  $\beta$ , proximal parameter  $\Delta\hat{\beta}$ , marginal error threshold  $\rho_{th}$ , target sparsity  $\lambda$ .

**Ensure:** The transport plan  $P$  for entropic OT with parameter  $\beta$

```

1:  $\hat{P}_\beta, u, v \leftarrow \text{IP-EOT}(\mathbf{a}, \mathbf{b}, C, \beta, \Delta\hat{\beta})$  ▷ # IP-EOT for approximation
2:  $\rho = \left\| \hat{P}_\beta \mathbf{1} - \mathbf{a} \right\|_1 + \left\| \hat{P}_\beta^\top \mathbf{1} - \mathbf{b} \right\|_1$ 
3: if  $\rho < \rho_{th}$  then  $P = \hat{P}_\beta$ 
4: else
5:    $\hat{H} \leftarrow \text{Sparsify Hessian}(\hat{P}_\beta, \rho)$  by (15) and Sec. 4.1
6:   if  $\|\hat{H}\|_0 < \lambda n$  then
7:      $P, \rho \leftarrow \text{Fast Sparse Newton}(\mathbf{a}, \mathbf{b}, C, \hat{P}_\beta, \rho_{th})$   $P, \rho \leftarrow \text{Fast Sparse Newton}(\mathbf{a}, \mathbf{b}, C, \hat{P}_\beta, \rho_{th}, \beta)$ 
8:   else
9:     while  $\rho > \rho_{th}$  do
10:       $\mathbf{u} \leftarrow \mathbf{a} \odot (\hat{P}_\beta \mathbf{v}), \mathbf{v} \leftarrow \mathbf{b} \odot (\hat{P}_\beta^\top \mathbf{u})$  ▷ # Sinkhorn scaling
11:      Compute  $\rho$  with  $P = \text{diag}(\mathbf{v}) \hat{P}_\beta \text{diag}(\mathbf{u})$ 
12:    end while
13:  end if
14: end if
15: return  $P$ 

```

---

**Datasets.** We consider three datasets following the experimental setup in (Tang et al., 2024): (1) Random linear assignment:  $n \in \{500, 4000\}$ , the cost matrix  $C \in \mathbb{R}^{n \times n}$  where are entries sampled from a uniform distribution over  $[0, 1]$ ; uniform marginals  $\mathbf{a} = \mathbf{b} = \frac{1}{500} \mathbf{1}$ ; This is a fundamental problem in graph matching. (2) MNIST: each image is flattened and normalized to form marginals; the cost matrix is defined by the squared Euclidean distance  $\|x - y\|_2$ ;  $n \in \{784, 4096\}$ ; (3) Multiple-solutions variant: based on MNIST, the cost matrix replaced with  $\ell_1$  norm  $\|x - y\|_1$  to remove uniqueness in the OT solution, which is a challenge for Newton-based algorithms (Tang et al., 2024). Finally, we set  $\beta = 200 \ln n$  for weak regularization, as the distance between  $P_\beta$  and the optimal plan  $P_\infty$  scales with  $\frac{\ln n}{\beta}$  (Altschuler et al., 2017). Under this setting, the values of  $\beta$  in experiments at the hundred-scale align with those reported in previous studies, making the comparison appropriate.

**Setting.** In PSN, the proximal parameter  $\Delta\hat{\beta}$  is set to 25 for general OT. This ensures numerical stability and leverages the rapid early convergence of IP-EOT, as the iteration count is  $\beta/\Delta\hat{\beta}$ . Given the faster convergence of the Sinkhorn method in linear assignment (Altschuler et al., 2017), we set  $\Delta\hat{\beta} = 50$  for this special case. The sparsity parameter  $\lambda$  is set to 30 for linear assignment and 70 for general OT. This design is informed by the theoretical sparsity of OT solutions ( $n$  for linear assignment and  $2n - 1$  for general OT) (Peyré et al., 2017) and the fact that entropic OT solutions are typically denser (Tang et al., 2024).

**Baselines** include the Sparse-Newton-Sinkhorn method (SNS, ICLR) (Tang et al., 2024), the Safe and Sparse Newton Method<sup>1</sup> (SSNS, NeurIPS) (Tang and Qiu, 2024), the Newton method for entropic OT<sup>2</sup> (Brauer et al., 2017), the log-Sinkhorn (Benamou et al., 2015), and the primal Sinkhorn (Cuturi, 2013) (when  $\beta > 100$ , the primal Sinkhorn may exhibit numerical instability, as reported in Appendix F of (Shen et al., 2024)). All methods were benchmarked in terms of runtime using MATLAB R2024b on an Intel Core i7-12800HX processor<sup>3</sup>.

---

<sup>1</sup><https://github.com/yixuan/regot-python>

<sup>2</sup><https://github.com/dirloren/sinkhornnewton>

<sup>3</sup>For SSNS implemented in C++, its runtime was scaled to the MATLAB baseline by comparing the execution time of the Newton method in C++ and MATLAB.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

### 6.1 IP-EOT VS. LOG-SINKHORN

This section compares IP-EOT with log-Sinkhorn. Figure 3 shows that the error of IP-EOT decreases much faster than that of log-Sinkhorn in the early iterations. Moreover, since the per-iteration cost of IP-EOT is less than that of log-Sinkhorn, its efficiency advantage in the early stage of the iterations becomes even more pronounced.

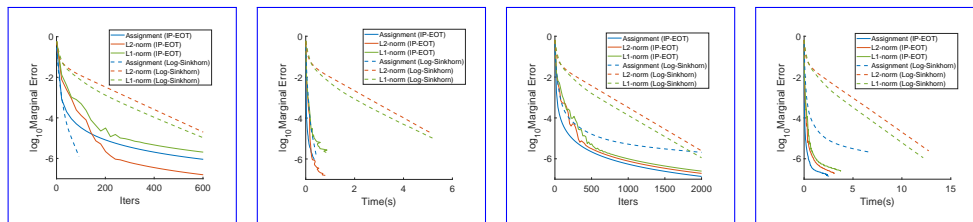


Figure 3: IP-EOT vs. log-Sinkhorn:  $\beta = 500$  for the first two subfigures,  $\beta = 1200$  for the latter two.

### 6.2 COMPARISON BETWEEN INCOMPLETE CHOLESKY AND DIAGONAL PRECONDITIONER

In this section, we compare this section compares the improvements provided by the Incomplete Cholesky (IC) and diagonal preconditioners when applied to the Newton method. We perform 300 Sinkhorn iterations across all cases to generate initial approximations. As shown in Figure 9, Figure 4 shows that the IC preconditioner enhances spectral clustering, significantly reduces CG iterations, and consistently outperforms the diagonal preconditioner across all datasets. Since Theorem 5 shows that CG’s per-step error is mainly governed by the clustered bulk of eigenvalues, these few isolated eigenvalues do not noticeably reduce the performance gains.

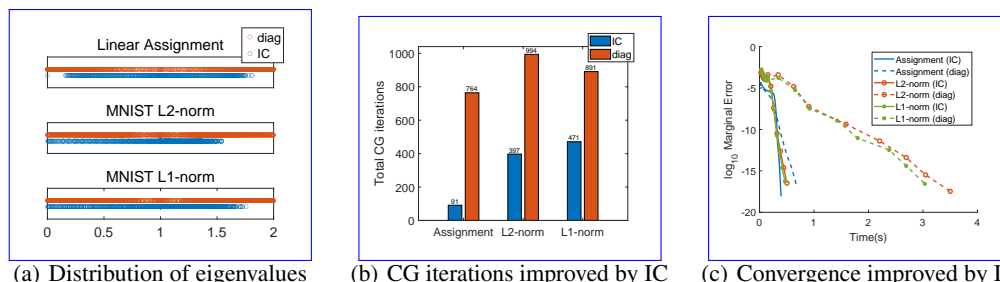


Figure 4: Comparison with diagonal and IC preconditioners. 300 Sinkhorn iterations are used to generate initial approximations.

### 6.3 COMPARISON BETWEEN PSN AND OTHER ALGORITHMS

This section evaluates PSN against several baseline methods. As illustrated in Figure 5, PSN demonstrates strong performance across all settings. In (a)–(c), PSN circumvents the numerical overflow risk inherent in the Sinkhorn method through a few proximal steps, achieving comparable efficiency safely and surpassing all other baselines. Subfigures (d)–(f) show that for large  $\beta$ , PSN outperforms other algorithms, especially log-Sinkhorn. In larger-scale problems (g)–(i), the Newton-based baselines suffer a sharp runtime increase, whereas PSN remains efficient and achieves over 20 $\times$  speedup over SNS and SSNS in some cases. Further results on iterations and time to accuracy are included in the Appendix.

## 7 CONCLUSION

We elaborate the trajectory of the regularization parameter during the iterations of the proximal point method and establish its ability to approximate the entropic optimal transport solution. This theoretical insight facilitates the seamless integration of the proximal point method with complementary algorithms, leading to the proposed Proximal-Sinkhorn-Newton framework that achieves high efficiency and robustness under varying entropic regularization parameters. Moreover, compared

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

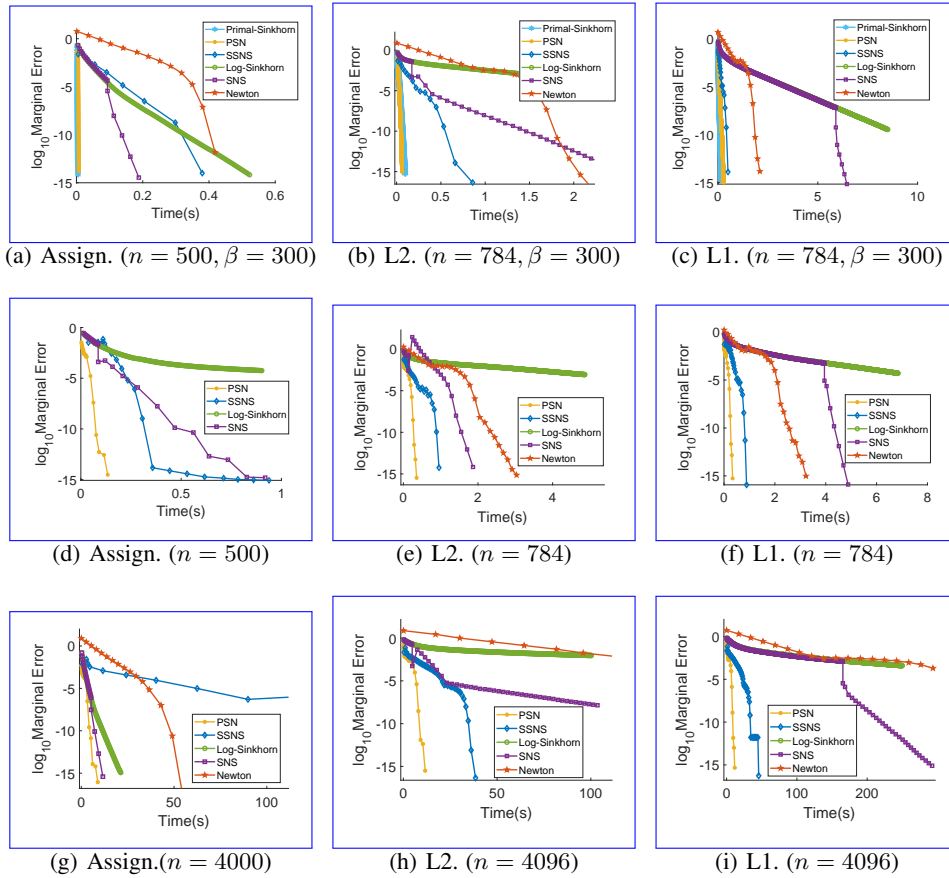


Figure 5: Comparisons between PSN and other algorithms on different problem sizes and different  $\beta$ . For subfigures without an explicit  $\beta$  value, the default setting  $\beta = 200 \ln(n)$  is adopted.

with previous Newton-based algorithms, the proposed fast Sparse Newton method accelerates each iteration with a preconditioner and mitigates stagnation during optimization through a novel line search strategy that incorporates Sinkhorn iterations scaling. Experimental results show that our approach PSN significantly faster than the state-of-the-art methods.

Since the number of IP-EOT iterations in this work is chosen empirically, future research may investigate how the number of proximal point iterations influences overall convergence, potentially leading to more adaptive iteration scheduling strategies. Another promising direction is the integration of additional algorithmic modules into the proximal point framework to further enhance performance. From a practical standpoint, exploiting hardware-based acceleration also holds considerable potential for improving computational efficiency.

---

594 **Ethics Statement.** This work focuses on the development and analysis of algorithms for entropic  
595 optimal transport. It does not involve human subjects, sensitive data, or applications with foreseeable  
596 ethical risks. To the best of our knowledge, this research raises no ethical concerns.  
597

598 **Reproducibility Statement.** To ensure reproducibility, we provide the source code with our sub-  
599 mission. All algorithmic steps, parameter settings, and experimental configurations are explicitly  
600 described in the paper. With the released code and detailed specifications, independent researchers  
601 should be able to fully reproduce our results.  
602

603 REFERENCES

604 Nathan Linial, Alex Samorodnitsky, and Avi Wigderson. A deterministic strongly polynomial  
605 algorithm for matrix scaling and approximate permanents. In Proceedings of the thirtieth annual  
606 ACM symposium on Theory of Computing, pages 644–652, 1998.  
607

608 Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport. Center for Research in  
609 Economics and Statistics Working Papers, 2017.

610 Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks.  
611 In International Conference on Machine Learning, pages 214–223. PMLR, 2017.  
612

613 Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn di-  
614 vergences. In International Conference on Artificial Intelligence and Statistics, pages 1608–1617.  
615 PMLR, 2018.

616 Binrui Shen, Qiang Niu, and Shengxin Zhu. Adaptive softassign via hadamard-equipped sinkhorn.  
617 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages  
618 17638–17647, 2024.

619 Tuan Nguyen, Trung Le, He Zhao, Quan Hung Tran, Truyen Nguyen, and Dinh Phung. Most:  
620 Multi-source domain adaptation via optimal transport for student-teacher learning. In Uncertainty  
621 in Artificial Intelligence, pages 225–235. PMLR, 2021.  
622

623 Rosanna Turrise, Rémi Flamary, Alain Rakotomamonjy, and Massimiliano Pontil. Multi-source  
624 Domain Adaptation via Weighted Joint Distributions Optimal Transport. In Uncertainty in Artificial  
625 Intelligence, pages 1970–1980. PMLR, 2022.

626 Luca Oneto, Michele Donini, Giulia Luise, Carlo Ciliberto, Andreas Maurer, and Massimiliano Pontil.  
627 Exploiting MMD and Sinkhorn Divergences for Fair and Transferable Representation Learning. In  
628 Advances in Neural Information Processing Systems, volume 33, pages 15360–15370, 2020.

629 Viet Huynh, He Zhao, and Dinh Phung. OTLDA: A Geometry-Aware Optimal Transport Approach  
630 for Topic Modeling. In Advances in Neural Information Processing Systems, volume 33, pages  
631 18573–18582, 2020.  
632

633 Na Lei, Kehua Su, Li Cui, Shing-Tung Yau, and Xianfeng David Gu. A Geometric View of Optimal  
634 Transportation and Generative Model. Computer Aided Geometric Design, 68:1–21, 2019.

635 Derek Onken, Samy Wu Fung, Xingjian Li, and Lars Ruthotto. OT-Flow: Fast and Accurate  
636 Continuous Normalizing Flows via Optimal Transport. In Proceedings of the AAAI Conference  
637 on Artificial Intelligence, 2021.  
638

639 Chunxu Cao and Qiang Zhang. Feature Selection via Maximizing Distances between Class Condi-  
640 tional Distributions. arXiv preprint arXiv:2401.07488, 2024.

641 Chunxu Cao, Yong Zhang, Yanke Ai, and Qiang Zhang. Unsupervised feature selection via unifying  
642 distribution alignment and structure preservation. Information Fusion, 126:103544, 2026. ISSN  
643 1566-2535. doi: <https://doi.org/10.1016/j.inffus.2025.103544>.

644 Gabriel Peyré and Marco Cuturi. Computational Optimal Transport: With Applications to Data  
645 Science. Foundations and Trends in Machine Learning, 2019.  
646

647 Marco Cuturi. Sinkhorn Distances: Lightspeed Computation of Optimal Transport. In  
Advances in Neural Information Processing Systems, volume 26, 2013.

---

648 Richard Sinkhorn. A Relationship Between Arbitrary Positive Matrices and Doubly Stochastic  
649 Matrices. The Annals of Mathematical Statistics, 35(2):876–879, 1964.

650

651 Jean-David Benamou, Guillaume Carlier, Marco Cuturi, Luca Nenna, and Gabriel Peyré. Iterative  
652 Bregman Projections for Regularized Transportation Problems. SIAM Journal on Scientific  
653 Computing, 37(2):A1111–A1138, 2015.

654 Yujia Xie, Xiangfeng Wang, Ruijia Wang, and Hongyuan Zha. A Fast Proximal Point Method  
655 for Computing Exact Wasserstein Distance. In Proceedings of The 35th Uncertainty in Artificial  
656 Intelligence Conference, volume 115, pages 433–453. PMLR, 22–25 Jul 2020.

657

658 Zihao Tang and Yixuan Qiu. Safe and sparse newton method for entropic-regularized optimal  
659 transport. In The Thirty-eighth Annual Conference on Neural Information Processing Systems,  
660 2024.

661 Christoph Brauer, Christian Clason, Dirk Lorenz, and Benedikt Wirth. A Sinkhorn-Newton method  
662 for entropic optimal transport. arXiv preprint arXiv:1710.06635, 2017.

663

664 Xun Tang, Michael Shavlovsky, Holakou Rahmanian, Elisa Tardini, Kiran Koshy Thekumparampil,  
665 Tesi Xiao, and Lexing Ying. Accelerating Sinkhorn algorithm with sparse Newton iterations. In  
666 The Twelfth International Conference on Learning Representations, 2024.

667 Sergey Guminov, Pavel Dvurechensky, Nazarii Tupitsa, and Alexander Gasnikov. On a Combination  
668 of Alternating Minimization and Nesterov’s Momentum. In International Conference on Machine  
669 Learning, pages 3886–3898. PMLR, 2021.

670 Jason Altschuler, Jonathan Niles-Weed, and Philippe Rigollet. Near-linear time approximation algo-  
671 rithms for optimal transport via sinkhorn iteration. In Advances in Neural Information Processing  
672 Systems, volume 30, 2017.

673

674 Pavel Dvurechensky, Alexander Gasnikov, Sergey Omelchenko, and Alexander Tiurin. A Stable Al-  
675 ternative to Sinkhorn’s Algorithm for Regularized Optimal Transport. In International Conference  
676 on Mathematical Optimization Theory and Operations Research, pages 406–423. Springer, 2020.

677 Weijian Li, Xianlin Zeng, and Lacro Pavel. Primal-dual Accelerated Mirror-Descent Method for  
678 Constrained Bilinear Saddle-Point Problems. IEEE Transactions on Automatic Control, 2025.

679

680 MJD Powell and Ph L Toint. On the estimation of sparse Hessian matrices. SIAM Journal on  
681 Numerical Analysis, 16(6):1060–1074, 1979.

682

683 Mete Kemertas, Amir-massoud Farahmand, and Allan Douglas Jepson. A Truncated Newton Method  
684 for Optimal Transport. In The Thirteenth International Conference on Learning Representations,  
685 2025.

686

687 Gene H Golub and Charles F Van Loan. Matrix Computations. JHU press, 2013.

688

689 Elena Caraba. Preconditioned Conjugate Gradient Algorithm, March 2008.

690

691 R Tyrrell Rockafellar. Monotone Operators and the Proximal Point Algorithm. SIAM Journal on  
692 Control and Optimization, 14(5):877–898, 1976.

693

694 Gong Chen and Marc Teboulle. Convergence Analysis of a Proximal-Like Minimization Algorithm  
695 Using Bregman Functions. SIAM Journal on Optimization, 3(3):538–543, 1993.

696

697 Hong TM Chu, Ling Liang, Kim-Chuan Toh, and Lei Yang. An efficient implementable inexact  
698 entropic proximal point algorithm for a class of linear programming problems. Computational  
699 Optimization and Applications, 85(1):107–146, 2023.

700

701 Lénaïc Chizat. Annealed Sinkhorn for Optimal Transport: convergence, regularization path and  
debiasing. arXiv preprint arXiv:2408.11620, 2024.

Jorge Nocedal and Stephen J Wright. Numerical Optimization. Springer, 1999.

Magnus R Hestenes, Eduard Stiefel, et al. Methods of Conjugate Gradients for Solving Linear  
Systems. Journal of research of the National Bureau of Standards, 49(6):409–436, 1952.

702 Yousef Saad. *Iterative Methods for Sparse Linear Systems*. SIAM, 2003.  
703  
704 Dong C Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization.  
705 *Mathematical programming*, 45(1):503–528, 1989.  
706 Pavel Dvurechensky, Alexander Gasnikov, and Alexey Kroshnin. Computational Optimal Trans-  
707 port: Complexity by Accelerated Gradient Descent Is Better Than by Sinkhorn’s Algorithm. In  
708 *International Conference on Machine Learning*, pages 1367–1376. PMLR, 2018.  
709 Zhi-Quan Luo and Paul Tseng. On the convergence of the coordinate descent method for convex  
710 differentiable minimization. *Journal of Optimization Theory and Applications*, 72(1):7–35, 1992.  
711

## 713 A NOTATIONS

714 Key notations are summarized below.

715 Table 1: Symbols and Notations.

| Symbol                                                                    | Definition                                                                                                                                |
|---------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------|
| $\mathcal{H}(\cdot)$                                                      | entropy function                                                                                                                          |
| $\mathcal{B}(\cdot)$                                                      | Bregman divergence                                                                                                                        |
| $\mathbb{U}$                                                              | $\{P \in \mathbb{R}_+^{n \times n} : P\mathbf{1} = \mathbf{a}, P^\top \mathbf{1} = \mathbf{b}\}, \mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ |
| $D_{\ell(\mathbf{x})}$                                                    | diagonal matrix of a vector $\mathbf{x}$                                                                                                  |
| $A \cdot B$                                                               | Frobenius inner product between matrix $A$ and $B$                                                                                        |
| $\otimes$                                                                 | Kronecker product                                                                                                                         |
| $\overset{\circ}{\exp}$                                                   | element-wise exponential                                                                                                                  |
| $\oslash$                                                                 | element-wise division                                                                                                                     |
| $\odot$                                                                   | element-wise product                                                                                                                      |
| $\bigcirc$                                                                | element-wise multiplication of a series of matrices                                                                                       |
| $\beta$                                                                   | the parameter in entropic OT                                                                                                              |
| $\Delta\beta$                                                             | the proximal parameter in proximal point methods                                                                                          |
| $\Gamma_{\mathbb{U}}^{\text{kl}}(\cdot), \hat{\Gamma}_{\text{kl}}(\cdot)$ | KL-projection and its approximating version                                                                                               |
| $P_\beta, \hat{P}_\beta$                                                  | the exact and approximating solution of OT with entropic parameter $\beta$                                                                |

## 731 B PROOFS FOR PROXIMAL POINT METHODS

732  
733 **Lemma 1 (Explicit form of  $\beta$ )** *The iterative schemes of the exact/inexact proximal point method*  
734 *for entropic OT admit the following explicit representation:*  
735

$$736 P^{(t)} = \Gamma_{\mathbb{U}}^{\text{kl}}\left(\bigcirc_{i=1}^t \overset{\circ}{\exp}(-C\beta^{(i)})\right), \quad \hat{P}^{(t)} = \hat{\Gamma}_{\mathbb{U}}^{\text{kl}}\left(\bigcirc_{i=1}^t \overset{\circ}{\exp}(-C\beta^{(i)})\right). \quad (24)$$

738 **Proof** *This proof begins with the iterative formula of the exact proximal point method:*

$$739 P^{(t)} = \Gamma_{\mathbb{U}}^{\text{kl}}\left(\overset{\circ}{\exp}(-C\beta^{(t)}) \odot P^{(t-1)}\right). \quad (25)$$

740 *For the first step, we have*

$$741 P^{(1)} = \Gamma_{\mathbb{U}}^{\text{kl}}(\overset{\circ}{\exp}(-C\beta^{(1)})) = D_{(\mathbf{u}_1)} \overset{\circ}{\exp}(-\beta^{(1)}C) D_{(\mathbf{v}_1)}, \quad (26)$$

742 *An equivalent form (Shen et al., 2024, Proposition 3) can be derived via elementary matrix operations*  
743 *as*

$$744 P^{(1)} = \overset{\circ}{\exp}(-C\beta^{(1)}) \odot (\mathbf{u}_1 \otimes \mathbf{v}_1^\top). \quad (27)$$

745 *Hence, for the subsequent iteration, we have:*

$$746 P^{(2)} = \Gamma_{\mathbb{U}}^{\text{kl}}\left(\overset{\circ}{\exp}(-C\beta^{(2)}) \odot P^{(1)}\right) \quad (28)$$

$$747 = \Gamma_{\mathbb{U}}^{\text{kl}}\left(\overset{\circ}{\exp}(-C\beta^{(2)}) \odot \overset{\circ}{\exp}(-C\beta^{(1)}) \odot (\mathbf{u}_1 \otimes \mathbf{v}_1^\top)\right) \quad (29)$$

$$748 = \left(\overset{\circ}{\exp}(-C\beta^{(2)}) \odot \overset{\circ}{\exp}(-C\beta^{(1)}) \odot (\mathbf{u}_1 \otimes \mathbf{v}_1^\top)\right) \odot (\mathbf{u}_2 \otimes \mathbf{v}_2^\top) \quad (30)$$

$$749 = \overset{\circ}{\exp}(-C(\beta^{(2)} + \beta^{(1)})) \odot (\mathbf{u}_1 \otimes \mathbf{v}_1^\top) \odot (\mathbf{u}_2 \otimes \mathbf{v}_2^\top) \quad (31)$$

$$750 = \overset{\circ}{\exp}(-C(\beta^{(2)} + \beta^{(1)})) \odot ((\mathbf{u}_1 \odot \mathbf{u}_2) \otimes (\mathbf{v}_1^\top \odot \mathbf{v}_2^\top)). \quad (32)$$

756 By defining  $\mathbf{u}_3 := \mathbf{u}_1 \odot \mathbf{u}_2$  and  $\mathbf{v}_3 := \mathbf{v}_1 \odot \mathbf{v}_2$ , we have

$$757 \quad P^{(2)} = \mathring{\text{exp}}(-C(\beta^{(2)} + \beta^{(1)})) \odot (\mathbf{u}_3 \otimes \mathbf{v}_3^\top). \quad (33)$$

759 Due to  $P^{(2)} \in \mathbb{U}$ ,  $P^{(2)}$  is the KL-projection of the kernel matrix  $\mathring{\text{exp}}(-C(\beta^{(1)} + \beta^{(2)}))$ . Since the  
760 KL-projection is unique for a given kernel matrix (Cuturi, 2013), it follows that

$$762 \quad P^{(2)} = \Gamma_{\mathbb{U}}^{\text{kl}}\left(\bigodot_{i=1}^2 \mathring{\text{exp}}(-C\beta^{(i)})\right). \quad (34)$$

764 The same structure holds for general  $t$  by induction:

$$766 \quad P^{(t)} = \Gamma_{\mathbb{U}}^{\text{kl}}\left(\bigodot_{i=1}^t \mathring{\text{exp}}(-C\beta^{(i)})\right). \quad (35)$$

769 For the inexact proximal point method, we can also obtain

$$771 \quad \hat{P}^{(2)} = \hat{\Gamma}_{\mathbb{U}}^{\text{kl}}\left(\mathring{\text{exp}}(-C\beta^{(2)}) \odot \hat{P}^{(1)}\right) \quad (36)$$

$$772 \quad = \hat{\Gamma}_{\mathbb{U}}^{\text{kl}}\left(\mathring{\text{exp}}(-C\beta^{(2)}) \odot \mathring{\text{exp}}(-C\beta^{(1)}) \odot (\hat{\mathbf{u}}_1 \otimes \hat{\mathbf{v}}_1^\top)\right) \quad (37)$$

$$774 \quad = \left(\mathring{\text{exp}}(-C\beta^{(2)}) \odot \mathring{\text{exp}}(-C\beta^{(1)}) \odot (\hat{\mathbf{u}}_1 \otimes \hat{\mathbf{v}}_1^\top)\right) \odot (\hat{\mathbf{u}}_2 \otimes \hat{\mathbf{v}}_2^\top) \quad (38)$$

$$776 \quad = \mathring{\text{exp}}(-C(\beta^{(2)} + \beta^{(1)})) \odot (\hat{\mathbf{u}}_1 \otimes \hat{\mathbf{v}}_1^\top) \odot (\hat{\mathbf{u}}_2 \otimes \hat{\mathbf{v}}_2^\top) \quad (39)$$

$$778 \quad = \mathring{\text{exp}}(-C(\beta^{(2)} + \beta^{(1)})) \odot \underbrace{(\hat{\mathbf{u}}_1 \odot \hat{\mathbf{u}}_2)}_{\hat{\mathbf{u}}_3} \otimes \underbrace{(\hat{\mathbf{v}}_1^\top \odot \hat{\mathbf{v}}_2^\top)}_{\hat{\mathbf{v}}_3^\top} \quad (40)$$

$$780 \quad = \mathring{\text{exp}}(-C(\beta^{(2)} + \beta^{(1)})) \odot (\hat{\mathbf{u}}_3 \otimes \hat{\mathbf{v}}_3^\top). \quad (41)$$

782 Since  $\hat{\mathbf{u}}_3$  and  $\hat{\mathbf{v}}_3$  are approximations for the scaling vectors of the corresponding kernel matrix  
783  $\mathring{\text{exp}}(-C(\beta^{(2)} + \beta^{(1)}))$ ,  $\hat{P}^{(2)}$  can be rewritten as

$$784 \quad \hat{P}^{(2)} = \hat{\Gamma}_{\mathbb{U}}^{\text{kl}}\left(\bigodot_{i=1}^2 \mathring{\text{exp}}(-C\beta^{(i)})\right). \quad (42)$$

788 The same structure holds for general  $t$  by induction:

$$790 \quad \hat{P}^{(t)} = \hat{\Gamma}_{\mathbb{U}}^{\text{kl}}\left(\bigodot_{i=1}^t \mathring{\text{exp}}(-C\beta^{(i)})\right). \quad (43)$$

793 **Lemma 2 (Inexact to Exact)** Let  $\hat{P}_\beta$  be an approximation to  $P_\beta$ , ~~then~~

$$795 \quad P_\beta = \Gamma_{\mathbb{U}}^{\text{kl}}(\hat{P}_\beta). \quad (44)$$

796 **Proof** Since  $\hat{P}_\beta$  is a approximation of  $P_\beta$ , we have

$$798 \quad \hat{P}_\beta = D_{(\hat{\mathbf{u}})} \mathring{\text{exp}}(-\beta C) D_{(\hat{\mathbf{v}})}, \quad (45)$$

799 which can be rewritten as

$$800 \quad \hat{P}_\beta = \mathring{\text{exp}}(-\beta C) \odot (\hat{\mathbf{u}} \otimes \hat{\mathbf{v}}^\top), \quad (46)$$

801 according to (Shen et al., 2024, Proposition 3). Applying the KL-projection yields

$$802 \quad \Gamma_{\mathbb{U}}^{\text{kl}}(\hat{P}_\beta) = \mathring{\text{exp}}(-\beta C) \odot (\hat{\mathbf{u}} \otimes \hat{\mathbf{v}}^\top) \odot (\mathbf{u}_2 \otimes \mathbf{v}_2^\top) \quad (47)$$

$$804 \quad = \mathring{\text{exp}}(-\beta C) \odot \underbrace{(\hat{\mathbf{u}} \odot \mathbf{u}_1)}_{\mathbf{u}_2} \otimes \underbrace{(\hat{\mathbf{v}}_1^\top \odot \hat{\mathbf{v}}^\top)}_{\mathbf{v}_2^\top} \quad (48)$$

$$806 \quad = \mathring{\text{exp}}(-\beta C) \odot (\mathbf{u}_2 \otimes \mathbf{v}_2^\top) \quad (49)$$

808 Since the KL-projection is unique for a given kernel matrix (Cuturi, 2013), it follows that  $\Gamma_{\mathbb{U}}^{\text{kl}}(\hat{P}_\beta) =$   
809  $\Gamma_{\mathbb{U}}^{\text{kl}}(\mathring{\text{exp}}(-\beta C)) = P_\beta$ .

810 **Theorem 1** *Inexact-to-Exact Transition*

811 *Let  $\hat{P}^{(t)}$  denotes the transport plan updated at the  $t$ -th iteration in the inexact proximal point*  
 812 *method with a fixed proximal parameter  $\Delta\beta$ , then*

$$813 P_\beta = \Gamma_{\mathbb{U}}^{\text{kl}}(\hat{P}^{(t)}), \text{ where } \beta = t\Delta\beta. \quad (50)$$

815 **Proof** *According to Lemma 1, we have*

$$816 \hat{P}^{(t)} = \hat{\Gamma}_{\text{kl}}(\text{e}\hat{\text{x}}\text{p}(-t\Delta\beta C)) = \hat{P}_\beta. \quad (51)$$

818 *Then, by Lemma 2, it follows that*

$$819 \Gamma_{\mathbb{U}}^{\text{kl}}(\hat{P}_\beta) = P_\beta, \quad (52)$$

820 *which completes the proof.*

824 **C PROOF FOR LOCAL QUADRATIC CONVERGENCE**

825 **Theorem 2 (Local Quadratic Convergence)** *Let the sequence  $\{z_k\}$  be generated by*

$$826 z_{k+1} = z_k - \hat{H}_k^{-1} \nabla f(z_k), \quad (53)$$

827 *where  $\hat{H}_k$  is sparsified from the Hessian matrix  $H_k = H(z_k)$  by the Fast-Definite-Positive Sparsify*  
 828 *Scheme. The sequence  $\{z_k\}$  converges locally quadratically to the minimizer  $z^*$ .*

829 **Proof** *Define that  $\|x\| = \sum_i |x_i|$  if  $x$  is a vector and  $\|X\| = \sum_i \sum_j |x_{ij}|$  if  $X$  is a matrix.*

830 *To establish quadratic convergence, we must show that there exists a constant  $c > 0$  such that for*  
 831 *sufficiently large  $k$ , the following inequality holds:*

$$832 \|z_{k+1} - z^*\| \leq c \|z_k - z^*\|^2. \quad (54)$$

833 *We begin the convergence analysis by examining the error at iteration  $k + 1$  and applying the triangle*  
 834 *inequality:*

$$835 \begin{aligned} 836 \|z_{k+1} - z^*\| &= \|z_k + \Delta z - z^*\| \\ 837 &= \|z_k - \hat{H}_k^{-1} \nabla f(z_k) - z^*\| \\ 838 &\leq \|\hat{H}_k^{-1}\| \|\hat{H}_k(z_k - z^*) - \nabla f(z_k)\| \\ 839 &\leq \|\hat{H}_k^{-1}\| \|\nabla f(z_k) - \nabla f(z^*) - H(z_k)(z_k - z^*) + H(z_k)(z_k - z^*) - \hat{H}_k(z_k - z^*)\| \\ 840 &\leq \underbrace{\|\hat{H}_k^{-1}\|}_{(I)} \left( \underbrace{\|\nabla f(z_k) - \nabla f(z^*) - H(z_k)(z_k - z^*)\|}_{(II)} + \underbrace{\|H(z_k) - \hat{H}_k\|}_{(III)} \|z_k - z^*\| \right). \end{aligned} \quad (55)$$

841 *If we can prove that exist positive constants  $c_1, c_2, c_3$ , such that*

$$842 \|\hat{H}_k^{-1}\| \leq c_1, \quad (56)$$

$$843 \|\nabla f(z_k) - \nabla f(z^*) - H(z_k)(z_k - z^*)\| \leq c_2 \|z_k - z^*\|^2, \quad (57)$$

$$844 \|H(z_k) - \hat{H}_k\| \leq c_3 \|z_k - z^*\|, \quad (58)$$

845 *then we get*

$$846 \|z_{k+1} - z^*\| \leq c_1 (c_2 + c_3) \|z_k - z^*\|^2, \quad (59)$$

847 *which implies quadratic convergence of  $\{z_k\}$ .*

848 *For the first part (I), we know that for symmetric matrix  $\hat{H}_k$ ,  $\|\hat{H}_k^{-1}\|_2 = 1/\lambda_{\min}(\hat{H}_k)$ . Therefore,*  
 849  *$\exists a > 0$ , we have*

$$850 \|\hat{H}_k^{-1}\| \leq \frac{a}{\lambda_{\min}(\hat{H}_k)}. \quad (60)$$

864 Now we use Gershgorin's Circle Theorem to estimate the lower bound of  $\lambda_{\min}(\hat{H}_k)$ . Define  $\delta_i =$   
865  $|\hat{h}_{ii}| - \sum_{j \neq i} |\hat{h}_{ij}|$ . Since the sparsification scheme eliminates the smallest entry at each step while  
866 preserving the diagonal elements, there exists  $r > 0$  such that  
867

$$868 \quad r < \delta_i = |h_{ii}| - \sum_{j \neq i} |\hat{h}_{ij}| = \left( \sum_{j \neq i} |h_{ij}| \right) - \left( \sum_{j \neq i} |\hat{h}_{ij}| \right) = \sum_{j \neq i} (|h_{ij}| - |\hat{h}_{ij}|). \quad (61)$$

872 Therefore, it follows directly from Gershgorin's Circle Theorem that

$$873 \quad \lambda_{\min}(\hat{H}_k) \geq \min_i \delta_i. \quad (62)$$

875 Let  $c_1 = a/r$ , then we get the inequality

$$877 \quad \|\hat{H}_k^{-1}\| \leq c_1. \quad (63)$$

879 For the second part (II), we define  $U(\mathbf{z}^*)$  is the neighborhood of  $\mathbf{z}^*$  and use the fact that  $\nabla f$  is  
880 continuously differentiable on  $U(\mathbf{z}^*)$ , and  $H(\mathbf{z})$  is Lipschitz continuous on the same neighborhood. This  
881 Lipschitz continuity implies that there exists a constant  $l$  such that

$$883 \quad \|H(\mathbf{z}_x) - H(\mathbf{z}_y)\| \leq l \|\mathbf{z}_x - \mathbf{z}_y\|, \quad \forall \mathbf{z}_x, \mathbf{z}_y \in U(\mathbf{z}^*). \quad (64)$$

884 Based on these properties, an application of Theorem 3.2.12 in [ortega2000iterative] shows that for  
885 any  $\mathbf{z}_k \in U(\mathbf{z}^*)$ , the following inequality holds:

$$887 \quad \|\nabla f(\mathbf{z}_k) - \nabla f(\mathbf{z}^*) - H(\mathbf{z}_k)(\mathbf{z}^* - \mathbf{z}_k)\| \leq \frac{l}{2} \|\mathbf{z}_k - \mathbf{z}^*\|^2. \quad (65)$$

889 Therefore, set  $c_2 = l/2$  and we get

$$891 \quad \|\nabla f(\mathbf{z}_k) - \nabla f(\mathbf{z}^*) - H(\mathbf{z}_k)(\mathbf{z}^* - \mathbf{z}_k)\| \leq c_2 \|\mathbf{z}_k - \mathbf{z}^*\|^2. \quad (66)$$

893 For the third part (III), we bound the error of the Hessian approximation,  $\|H_k - \hat{H}_k\|$ . Let  $E_k =$   
894  $H_k - \hat{H}_k$ , and let  $e_{ij}^{(k)}$  denote an entry of  $E_k$ . We divide  $e_{ij}^{(k)}$  into two parts: one corresponding to  
895 the finite set  $\mathcal{S}$  of index pairs  $(i, j)$  where  $h_{ij}^* = 0$ , and the other corresponding to the set  $\bar{\mathcal{S}}$ , which  
896 contains index pairs where  $h_{ij}^* \neq 0$ . Therefore, we have

$$898 \quad \|H_k - \hat{H}_k\| = \sum_i \sum_j |e_{ij}^{(k)}| = \sum_{(i,j) \in \mathcal{S}} |e_{ij}^{(k)}| + \sum_{(i,j) \in \bar{\mathcal{S}}} |e_{ij}^{(k)}|. \quad (67)$$

902 When  $(i, j) \in \mathcal{S}$ , we have  $h_{ij}^* = 0$ . We use Lipschitz continuity of  $H(\mathbf{z})$  to bound the term  $|e_{ij}^{(k)}|$ ,

$$904 \quad |e_{ij}^{(k)}| = |h_{ij}^{(k)} - h_{ij}^*| \leq \|H_k - H(\mathbf{z}^*)\| \leq l \|\mathbf{z}_k - \mathbf{z}^*\|. \quad (68)$$

906 Let the size of  $\mathcal{S}$  be  $d$ . We can bound the norm of  $E_k$  by defining  $c_3 = dl > 0$ ,

$$908 \quad \sum_{(i,j) \in \mathcal{S}} |e_{ij}^{(k)}| \leq d \cdot l \|\mathbf{z}_k - \mathbf{z}^*\| = c_3 \|\mathbf{z}_k - \mathbf{z}^*\|. \quad (69)$$

910 For the second part of (67), we have

$$912 \quad \sum_{(i,j) \in \bar{\mathcal{S}}} |e_{ij}^{(k)}| \leq \sum_{(i,j) \in \bar{\mathcal{S}}} \tau = c_4 \tau. \quad (70)$$

915 where  $c_4 = 2n^2 - d$ . Combining (69) and (70) we get that

$$917 \quad \sum_i \sum_j |e_{ij}^{(k)}| \leq c_3 \|\mathbf{z}_k - \mathbf{z}^*\| + c_4 \tau. \quad (71)$$

918 Since  $\tau = \|\mathbf{z}_k - \mathbf{z}^*\|$ , it follows that

$$919 \sum_i \sum_j \left| e_{ij}^{(k)} \right| \leq (c_3 + c_4) \|\mathbf{z}_k - \mathbf{z}^*\|, \quad (72)$$

922 so we have

$$923 \|H_k - \hat{H}_k\| = \sum_i \sum_j \left| e_{ij}^{(k)} \right| \leq (c_3 + c_4) \|\mathbf{z}_k - \mathbf{z}^*\|. \quad (73)$$

926 Combining (63), (66), (73) and the formula of  $\|\mathbf{z}_{k+1} - \mathbf{z}^*\|$ , we obtain

$$928 \|\mathbf{z}_{k+1} - \mathbf{z}^*\| \leq c_1 (c_2 + c_3 + c_4) \|\mathbf{z}_k - \mathbf{z}^*\|^2, \quad (74)$$

929 implying the local quadratic convergence.

## 932 D PROOF FOR GLOBAL CONVERGENCE

933 **Lemma 3** Let  $f(\mathbf{z})$  be defined in equation (14), and the sequence of search directions  $\{\mathbf{p}_k\}$  be generated by the Sinkhorn iterations. Then in each update, direction  $\mathbf{p}_k$  with unit step size  $\alpha_k = 1$  satisfy satisfies the following conditions:

$$937 f(\mathbf{z}_k + \mathbf{p}_k) \leq f(\mathbf{z}_k) + c_1 \nabla f(\mathbf{z}_k)^\top \mathbf{p}_k, \quad (75)$$

$$939 \nabla f(\mathbf{z}_k + \mathbf{p}_k)^\top \mathbf{p}_k \geq c_2 \nabla f(\mathbf{z}_k)^\top \mathbf{p}_k, \quad (76)$$

940 for constants  $0 < c_1 < c_2 < 1$ .

942 **Proof** Since Sinkhorn algorithm is a Block Coordinate Descent method, equation (75) is naturally satisfied. Let us analyze a single step of updating the  $\mathbf{y}$  block while keeping the  $\mathbf{x}$  block fixed, and symmetric argument applies to the update of the  $\mathbf{x}$  block.

944 A Sinkhorn iteration updates the  $\mathbf{y}$  block by solving the subproblem:

$$947 \mathbf{y}_{k+1} = \arg \min_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{x}_{k+1}, \mathbf{y}). \quad (77)$$

949 The first-order optimality condition for this subproblem requires the partial gradient with respect to  $\mathbf{y}$  is zero, i.e.,

$$951 \nabla_{\mathbf{y}} f(\mathbf{x}_{k+1}, \mathbf{y}_{k+1}) = 0. \quad (78)$$

952 The search direction  $\mathbf{p}_k$  associated with this update has a block structure  $\mathbf{p}_k = (0, \mathbf{p}_k^{\mathbf{y}})$ , as it only modifies the  $\mathbf{y}$  components. Consequently, the projection of the new gradient  $\nabla f(\mathbf{z}_{k+1})$  onto the search direction  $\mathbf{p}_k$  is:

$$955 \nabla f(\mathbf{z}_{k+1})^\top \mathbf{p}_k = \nabla_{\mathbf{x}} f(\mathbf{z}_{k+1})^\top \cdot 0 + \nabla_{\mathbf{y}} f(\mathbf{z}_{k+1})^\top \mathbf{p}_k^{\mathbf{y}} = 0. \quad (79)$$

957 Since the Sinkhorn iteration employs a descent direction, we have  $\nabla f(\mathbf{z}_k)^\top \mathbf{p}_k < 0$ . Combining with equation (79) it leads to

$$959 c_2 \nabla f(\mathbf{z}_k)^\top \mathbf{p}_k < 0 = \nabla f(\mathbf{z}_{k+1})^\top \mathbf{p}_k, \quad (80)$$

960 which satisfies the equation (76) for any  $c_2 \in (c_1, 1)$ . Thus, the conditions of this Lemma are always met by Sinkhorn iterations.

963 **Lemma 4** Let  $f(\mathbf{z})$  be defined in equation (14), the sequences  $\{\mathbf{z}_k\}$  and  $\{\alpha_k\}$  be generated by Algorithm 2, and  $\mathbf{p}_k$  be the search direction. Then, the function  $f(\mathbf{z})$  and the generated sequences possess the following properties:

- 966 •  $f(\mathbf{z})$  is bounded below, i.e., there exists a scalar  $\bar{f}$  such that  $f(\mathbf{z}) \geq \bar{f}$  for all  $\mathbf{z}$ .
- 967 •  $f(\mathbf{z})$  is continuously differentiable.
- 968 • The gradient of  $f$ ,  $\nabla f$ , is globally Lipschitz continuous, i.e., there exists a constant  $l > 0$  such that

$$970 \|\nabla f(\mathbf{z}_x) - \nabla f(\mathbf{z}_y)\| \leq l \|\mathbf{z}_x - \mathbf{z}_y\|, \quad \forall \mathbf{z}_x, \mathbf{z}_y. \quad (81)$$

- 972 • The step size  $\alpha_k$  is determined by a line search procedure that ensures the Wolfe conditions  
 973 are satisfied for all  $k$ :

974  
 975 
$$f(\mathbf{z}_k + \alpha_k \mathbf{p}_k) \leq f(\mathbf{z}_k) + c_1 \alpha_k \nabla f(\mathbf{z}_k)^\top \mathbf{p}_k, \quad (82)$$

976 
$$\nabla f(\mathbf{z}_k + \alpha_k \mathbf{p}_k)^\top \mathbf{p}_k \geq c_2 \nabla f(\mathbf{z}_k)^\top \mathbf{p}_k, \quad (83)$$

977 for constants  $0 < c_1 < c_2 < 1$ . Specially,  $\mathbf{p}_k = \Delta \mathbf{z}_k$  when line search method satisfies  
 978 Wolfe condition.  
 979

980 Consequently, the sequence of gradients  $g_k = \nabla f(\mathbf{z}_k)$  converges to zero, that is,

981  
 982 
$$\lim_{k \rightarrow \infty} \|g_k\| = 0. \quad (84)$$

983 **Proof** By the properties of the equation (1), the first three conditions are already met. Algorithm 2  
 984 choose Sinkhorn iteration to decrease  $f(\mathbf{z})$  if equation (83) is not met by line search. But in Lemma  
 985 3 we have already proven that Sinkhorn iterations-scaling satisfy Wolfe condition, so Algorithm 2  
 986 possess these four conditions.

987 Given that  $f(\mathbf{z})$  is bounded below, continuously differentiable with a Lipschitz continuous gradient,  
 988 and the iterates satisfy the Wolfe conditions, we can apply Zoutendijk's Theorem, which states that:

989  
 990 
$$\sum_{k=0}^{\infty} \cos^2 \theta_k \|g_k\|^2 < \infty, \quad (85)$$

991 where  $\theta_k$  is the angle between the search direction  $\Delta \mathbf{z}_k$  and the negative gradient  $-g_k$ , defined by  
 992  $\cos \theta_k = \frac{-g_k^\top \Delta \mathbf{z}_k}{\|g_k\| \|\Delta \mathbf{z}_k\|}$ . The convergence  $\lim_{k \rightarrow \infty} \|g_k\| = 0$  follows if  $\cos \theta_k$  is uniformly bounded by  
 993 positive constant.

994 When line search method find the step size  $\alpha_k$  satisfies Wolfe condition, the search direction is given  
 995 by  $\Delta \mathbf{z}_k = -\hat{H}_k^{-1} g_k$ , where  $\hat{H}_k$  is a positive definite matrix. From Gershgorin's Circle Theorem we  
 996 know that the maximum eigenvalue of  $\hat{H}_k$  is bounded by  $2 \max_{i,j} (\mathbf{a}_i, \mathbf{b}_j)$ , so combine with equation  
 997 (63) we can get that exists constant  $c_{\text{line search}} > 0$  such that  $\kappa(\hat{H}_k) \leq c_{\text{line search}}$ .

1000 Then we establish a positive lower bound for  $\cos \theta_k$ :

1001  
 1002 
$$\cos \theta_k = \frac{-g_k^\top (-\hat{H}_k^{-1} g_k)}{\|g_k\| \|\hat{H}_k^{-1} g_k\|} = \frac{g_k^\top \hat{H}_k^{-1} g_k}{\|g_k\| \|\hat{H}_k^{-1} g_k\|}. \quad (86)$$

1003 Using the Rayleigh quotient and properties of matrix norms, we bound the numerator and denomina-  
 1004 tor:

- 1005 • Numerator:  $g_k^\top \hat{H}_k^{-1} g_k \geq \lambda_{\min}(\hat{H}_k^{-1}) \|g_k\|^2 = \frac{1}{\lambda_{\max}(\hat{H}_k)} \|g_k\|^2$ .  
 1006 • Denominator:  $\|\hat{H}_k^{-1} g_k\| \leq \|\hat{H}_k^{-1}\|_2 \|g_k\| = \lambda_{\max}(\hat{H}_k^{-1}) \|g_k\| = \frac{1}{\lambda_{\min}(\hat{H}_k)} \|g_k\|$ .

1007 Combining these inequalities yields:

1008  
 1009 
$$\cos \theta_k \geq \frac{\frac{1}{\lambda_{\max}(\hat{H}_k)} \|g_k\|^2}{\|g_k\| \cdot \frac{1}{\lambda_{\min}(\hat{H}_k)} \|g_k\|} = \frac{\lambda_{\min}(\hat{H}_k)}{\lambda_{\max}(\hat{H}_k)} = \frac{1}{\kappa(\hat{H}_k)}. \quad (87)$$

1010 When the line search fails, Algorithm 2 employs a Sinkhorn step as a fallback. Now we demonstrate  
 1011 that this step can also be framed as a preconditioned gradient descent, where the corresponding  
 1012 preconditioner has a uniformly bounded condition number.

1013 In the Sinkhorn iteration, the descent direction  $\mathbf{p}_k$  can be defined by:

1014  
 1015 
$$\mathbf{p}_k = -\tilde{H}_k^{-1} g_k, \quad (88)$$

1016 where  $\tilde{H}_k$  is the matrix given by:

1017  
 1018 
$$\tilde{H}_k = \beta \begin{pmatrix} \text{diag}(\mathcal{D}(P_k \mathbf{1}_m, \mathbf{a})) & \mathbf{0} \\ \mathbf{0} & \text{diag}(\mathcal{D}(P_k^\top \mathbf{1}_n, \mathbf{b})) \end{pmatrix}, \quad (89)$$

$\mathcal{D}(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y}) \odot (\log \mathbf{x} - \log \mathbf{y})$  and  $P_k$  is the transport plan at the  $k$ -th iteration.

Let us analyze the properties of  $\tilde{H}_k$ . Since the target marginals  $\mathbf{a}$  and  $\mathbf{b}$  consist of strictly positive elements, the row sums and column sums of the transport plan  $P_k$  also remain strictly positive throughout the iterations. Consequently,  $\tilde{H}_k$  is a diagonal matrix with strictly positive entries, which ensures it is positive definite.

The eigenvalues of  $\tilde{H}_k$  are its diagonal entries, specifically  $\{\beta \mathcal{D}(P_k \mathbf{1}_m, \mathbf{a})\}_i \cup \{\beta \mathcal{D}(P_k^\top \mathbf{1}_n, \mathbf{b})\}_j$ . Therefore, there exist constants  $\lambda_{\min}^* > 0$  and  $\lambda_{\max}^* < \infty$ , independent of  $k$ , such that all eigenvalues of  $\tilde{H}_k$  lie within the interval  $[\lambda_{\min}^*, \lambda_{\max}^*]$ . This implies that its condition number is uniformly bounded:

$$\kappa(\tilde{H}_k) = \frac{\lambda_{\max}(\tilde{H}_k)}{\lambda_{\min}(\tilde{H}_k)} \leq \frac{\lambda_{\max}^*}{\lambda_{\min}^*} =: c_{\text{Sinkhorn}}. \quad (90)$$

Substituting this result into the same inequality for  $\cos \theta_k$  as in the former case, we obtain:

$$\cos \theta_k \geq \frac{1}{\kappa(\tilde{H}_k)} \geq \frac{1}{c_{\text{Sinkhorn}}} > 0. \quad (91)$$

Since  $\kappa(\hat{H}_k) \leq c_{\text{line search}}$ ,  $\kappa(\tilde{H}_k) \leq c_{\text{Sinkhorn}}$ , let  $c = \max(c_{\text{line search}}, c_{\text{Sinkhorn}})$ , we have  $\cos \theta_k \geq \frac{1}{c} > 0$ , which provides a uniform positive lower bound. Applying this to Zoutendijk's condition:

$$\infty > \sum_{k=0}^{\infty} \cos^2 \theta_k \|g_k\|^2 \geq \sum_{k=0}^{\infty} \frac{1}{c^2} \|g_k\|^2 = \frac{1}{c^2} \sum_{k=0}^{\infty} \|g_k\|^2. \quad (92)$$

This implies that the series  $\sum_{k=0}^{\infty} \|g_k\|^2$  must converge. Therefore,

$$\lim_{k \rightarrow \infty} \|g_k\|^2 = 0, \quad (93)$$

which implies:

$$\lim_{k \rightarrow \infty} \|g_k\| = 0. \quad (94)$$

**Theorem 3 (Global Convergence)** The sequence  $\{\mathbf{z}_k\}$  generated by Algorithm 2 converges globally to the minimizer  $\mathbf{z}^*$ .

**Proof** Now we begin to prove global convergence. Let  $\mathbf{z}^*$  be the global optimum of (14). Since  $f$  and  $\nabla f$  are both Lip-continuous,  $\exists c_1, c_2 > 0$  such that

$$|f(\mathbf{z}_k) - f(\mathbf{z}^*)| \leq c_1 \|\mathbf{z}_k - \mathbf{z}^*\|, \quad (95)$$

$$\|\nabla f(\mathbf{z}_k) - \nabla f(\mathbf{z}^*)\| \leq c_2 \|\mathbf{z}_k - \mathbf{z}^*\|. \quad (96)$$

Firstly, Tang and Qiu (2024) show that the minimum eigenvalue of the Hessian matrix of the function  $f$  define by (14) has a positive lower bound, i.e.  $\exists c_3 > 0$ ,  $\lambda_{\min}(H(\mathbf{z})) \geq c_3$ , so we can get that

$$(\nabla f(\mathbf{z}_k) - \nabla f(\mathbf{z}^*))^\top (\mathbf{z}_k - \mathbf{z}^*) \geq c_3 \|\mathbf{z}_k - \mathbf{z}^*\|^2, \quad (97)$$

which implies

$$\|\nabla f(\mathbf{z}_k)\| = \|\nabla f(\mathbf{z}_k) - \nabla f(\mathbf{z}^*)\| \geq c_3 \|\mathbf{z}_k - \mathbf{z}^*\|. \quad (98)$$

On the other hand, by Taylor's theorem,

$$f(\mathbf{z}_k) = f(\mathbf{z}^*) + (\nabla f(\mathbf{z}_k))^\top (\mathbf{z}_k - \mathbf{z}^*) + \frac{1}{2} (\mathbf{z}_k - \mathbf{z}^*)^\top H(\tilde{\mathbf{z}}) (\mathbf{z}_k - \mathbf{z}^*), \quad (99)$$

where  $\tilde{\mathbf{z}}$  is a point between  $\mathbf{z}_k$  and  $\mathbf{z}^*$ , so we have

$$|f(\mathbf{z}_k) - f(\mathbf{z}^*)| \geq \frac{c_3}{2} \|\mathbf{z}_k - \mathbf{z}^*\|^2. \quad (100)$$

Since  $f$  is non-increase and has lower bound,  $\exists \epsilon \geq 0$  such that

$$\lim_{k \rightarrow \infty} |f(\mathbf{z}_k) - f(\mathbf{z}^*)| = \epsilon. \quad (101)$$

Combine (95) and (101), for sufficient large  $k$  we have

$$\frac{1}{2} c_1^{-1} \epsilon \leq c_1^{-1} |f(\mathbf{z}_k) - f(\mathbf{z}^*)| \leq \|\mathbf{z}_k - \mathbf{z}^*\| \leq c_3^{-1} \|\nabla f(\mathbf{z}_k)\|. \quad (102)$$

Combine (96), (100) and (101), for sufficient large  $k$  we have

$$\frac{c_3}{2c_2^2} \|\nabla f(\mathbf{z}_k)\|^2 \leq \frac{c_3}{2} \|\mathbf{z}_k - \mathbf{z}^*\|^2 \leq |f(\mathbf{z}_k) - f(\mathbf{z}^*)| \rightarrow \epsilon. \quad (103)$$

Since our line search method and Sinkhorn iteration both satisfy Wolfe condition, applying Lemma 4 we get that  $\epsilon = 0$ . Combine (102) and (103) with squeeze theorem, we have

$$\lim_{k \rightarrow \infty} \|\mathbf{z}_k - \mathbf{z}^*\| = 0. \quad (104)$$

**Theorem 4** [Algorithm 3 converges globally to the solution of the corresponding entropic OT.](#)

**Proof** [The global convergence of the PSN algorithm follows from the global convergence of its two components: the Proximal-Sinkhorn \(PS\) and Proximal-Newton \(PN\). The convergence of PN is established by Theorem 3, while that of PS is guaranteed by the global convergence of the Sinkhorn algorithm Cutri \(2013\). Hence, PSN is globally convergent.](#)

## E THE SPECTRAL DISTRIBUTION OF THE HESSIAN MATRIX

The following Lemma and Theorem characterize the convergence behavior of CG. Lemma 5 shows that the error at each iteration of CG is bounded by a polynomial  $q_k(\lambda)$ , and Theorem 5 reveals the connection between the polynomial  $q_k(\lambda)$  and the spectral distribution of the Hessian matrix  $H$  while solving  $H\Delta\mathbf{z} = -\nabla f$ . **Furthermore**, Theorem 5 also demonstrates that the error vector in each iteration of the CG method remains bounded by the polynomial  $q_k(\lambda)$ , even after excluding a small number of eigenvalues at one or both extremes. This reveals that the CG method's convergence rate is closely tied to the eigenvalue distribution of the matrix  $H$ , potentially leading to superlinear rates.

**Lemma 5 ((Saad, 2003))** *The approximate solution  $\mathbf{x}_k$  generated by the CG method of solving  $H\mathbf{x} = \mathbf{b}$  at step  $k$  satisfies*

$$\|\mathbf{x}_k - \mathbf{x}^*\|_H \leq \min_{q_k \in \mathcal{P}_k^{(0)}} \max_{\lambda \in \lambda(H)} |q_k(\lambda)| \|\mathbf{x}_0 - \mathbf{x}^*\|_H, \quad (105)$$

where  $\mathbf{x}^*$  denotes the exact solution of the linear system,  $\|\cdot\|_H$  represents the  $H$ -norm,  $\mathcal{P}_k^{(0)}$  represents the set of all real coefficient polynomials of degree not exceeding  $k$  that satisfy  $q_k(0) = 1$ .

**Theorem 5 ((Hestenes et al., 1952))** *Considering the error vectors  $\mathbf{u}_k = \mathbf{x}_k - \mathbf{x}^*$  defined by Lemma 5, we establish the following bounds:*

(1) for  $\lambda(H) \subset \{\lambda_1, \dots, \lambda_p\} \cup [a_1, b_1]$ , where  $\lambda_i \in \lambda(H)$ ,  $\lambda_i < a_1$ , we have

$$\|\mathbf{u}_k\|_H \leq F(a_1, b_1; k-p) \prod_{i=1}^p \frac{b_1}{\lambda_i} \|\mathbf{u}_0\|_H; \quad (106)$$

(2) for  $\lambda(H) \subset [a_2, b_2] \cup \{\lambda'_1, \dots, \lambda'_p\}$ , where  $\lambda'_i \in \lambda(H)$ ,  $\lambda'_i > b_2$ , we have

$$\|\mathbf{u}_k\|_H \leq F(a_2, b_2; k-p) \|\mathbf{u}_0\|_H; \quad (107)$$

(3) for  $\lambda(H) \subset \{\lambda_1, \dots, \lambda_p\} \cup [a_3, b_3] \cup \{\lambda'_1, \dots, \lambda'_p\}$ , where  $\lambda_i, \lambda'_i \in \lambda(H)$ ,  $\lambda'_i > b_3$ ,  $\lambda_i < a_3$ , we have

$$\|\mathbf{u}_k\|_H \leq \frac{1}{4} \prod_{i=1}^p \frac{\lambda'_i}{\lambda_i} \left(1 - \frac{\lambda_i}{\lambda'_i}\right)^2 F(a_3, b_3; k-2) \|\mathbf{u}_0\|_H, \quad (108)$$

where  $F(a, b; k) = \left(T_k\left(\frac{b+a}{b-a}\right)\right)^{-1}$ , and  $T_k(x)$  denotes Chebyshev polynomial of degree  $k$ .

## F PERFORMANCE ON DIFFERENT ENTROPIC PARAMETERS

We evaluate PSN, SNS, SSNS, and the Newton method across entropic parameters  $\beta \in \{100, 500, 3000\}$  covering both low and high regularization. Figure 6 indicates that our proposed algorithm, PSN, exhibits superior robustness and effectiveness when  $\beta$  is large. For smaller values of  $\beta$ , PSN reduces to a combination of the proximal point method and Sinkhorn [iterationscaling](#), and outperforms other Newton-based algorithms as well as the log-Sinkhorn method.

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

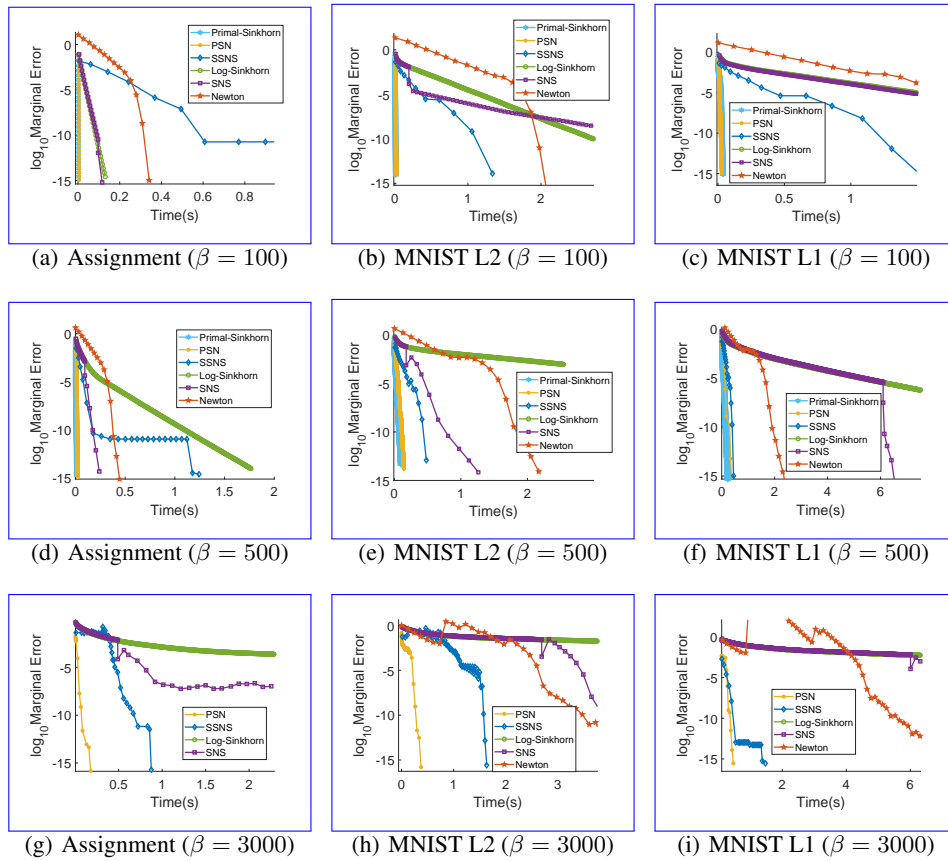


Figure 6: Performance on different  $\beta$ .

## G COMPARISON WITH OTHER CLASSIC FIRST ORDER ALGORITHMS

In this section, we compare our algorithm against several well-established first order optimization methods, including Limited-memory Broyden–Fletcher–Goldfarb–Shanno algorithm (L-BFGS, (Liu and Nocedal, 1989)) for the quasi-Newton method, the adaptive primal-dual accelerated gradient descent method (ADPAGD, (Dvurechensky et al., 2018)) as an accelerated first-order approach, and the Sinkhorn method implemented as a block coordinate descent method (BCD, (Luo and Tseng, 1992)). We adopted the C++ implementation by Tang and Qiu (2024) for the baseline algorithms<sup>4</sup>.

Evaluated across our three datasets, our algorithm, as depicted in Figure 7, consistently delivers superior performance, significantly outperforming these classic optimization techniques in all situations.

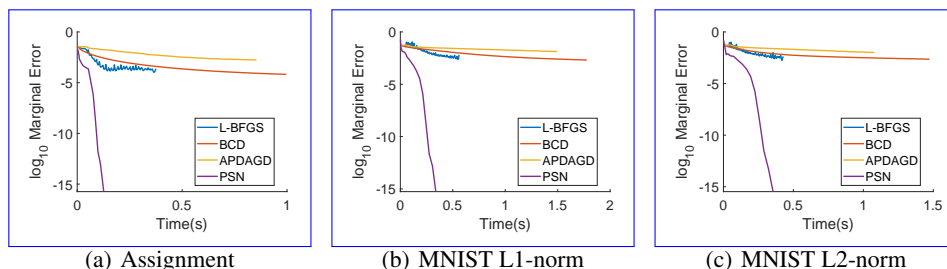


Figure 7: Other classic algorithms.

## H TRUNCATION DETAILS

We explore the sparsity of the Hessian matrix under various truncation thresholds and visualize the results in Figure 8.

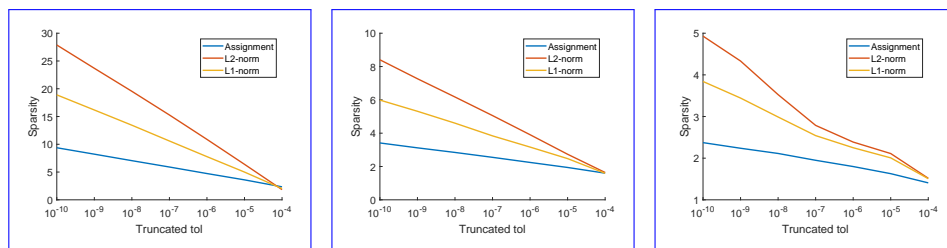


Figure 8: Relationship between truncation threshold and Sparsity. The [three subfigures correspond to  \$\beta = \{1000, 2000, 4000\}\$ , respectively](#). The  $x$ -axis denotes the truncation threshold, while the  $y$ -axis represents the Hessian matrix sparsity as a multiple of  $1/n$ .

Figure 8 demonstrates that the truncated matrix exhibits significant sparsity. ~~The number of non-zero elements after applying the  $10^{-4}$  threshold is comparable to that of the optimal OT solution, suggesting that a larger truncation threshold may lead to excessive sparsification.~~ [Moreover, under the same truncation threshold, larger values of  \$\beta\$  lead to a sparser matrix.](#)

[Interestingly, the relationship between the truncation threshold and the resulting sparsity level is nearly linear. Consequently, when the sparsified matrix contains fewer than  \$6n\$  nonzero entries, we can readily determine the truncation threshold that yields approximately  \$6n\$  nonzeros.](#)

<sup>4</sup><https://github.com/yixuan/regot-python>.

## I DETAILED RESULTS OF FIGURE 5

| Case                    | Error     | IP-EOT  |       | Primal-Sinkhorn |        | Log-Sinkhorn |       | PSN     |        | SSNS      |       | SNS                |       | Newton    |       | SSNS/PSN |
|-------------------------|-----------|---------|-------|-----------------|--------|--------------|-------|---------|--------|-----------|-------|--------------------|-------|-----------|-------|----------|
|                         |           | # iters | time  | # iters         | time   | # iters      | time  | # iters | time   | # iters   | time  | # iters            | time  | # iters   | time  |          |
| Assignment<br>$n = 500$ | $10^{-4}$ | 16      | 0.01s | 12              | 0.001s | 12           | 0.06s | 11      | 0.005s | <u>6</u>  | 0.14s | 12                 | 0.06s | <u>12</u> | 0.32s | 28×      |
|                         | $10^{-6}$ | 148     | 0.08s | 25              | 0.002s | 25           | 0.10s | 23      | 0.006s | <u>7</u>  | 0.20s | <u>20</u> <u>1</u> | 0.09s | <u>13</u> | 0.35s | 33×      |
|                         | $10^{-8}$ | 1466    | 0.76s | 43              | 0.003s | 43           | 0.19s | 42      | 0.007s | <u>8</u>  | 0.30s | <u>20</u> <u>2</u> | 0.11s | <u>14</u> | 0.38s | 42×      |
| MNIST L2<br>$n = 784$   | $10^{-4}$ | 39      | 0.04s | 176             | 0.02s  | 176          | 1.48s | 109     | 0.02s  | <u>10</u> | 0.22s | <u>20</u> <u>1</u> | 0.17s | <u>13</u> | 1.45s | 11×      |
|                         | $10^{-6}$ | 111     | 0.12s | 392             | 0.04s  | 392          | 3.33s | 324     | 0.03s  | <u>13</u> | 0.38s | <u>20</u> <u>4</u> | 0.4s  | <u>14</u> | 1.57s | 13×      |
|                         | $10^{-8}$ | 671     | 0.77s | 608             | 0.05s  | 608          | 5.16s | 540     | 0.04s  | <u>15</u> | 0.54s | <u>20</u> <u>9</u> | 0.82s | <u>15</u> | 1.69s | 14×      |
| MNIST L1<br>$n = 784$   | $10^{-4}$ | 57      | 0.06s | 148             | 0.02s  | 148          | 1.25s | 114     | 0.04s  | <u>10</u> | 0.16s | 148                | 1.26s | <u>13</u> | 1.51s | 4×       |
|                         | $10^{-6}$ | 296     | 0.32s | 409             | 0.05s  | 409          | 3.50s | 354     | 0.07s  | <u>15</u> | 0.36s | 409                | 3.45s | <u>15</u> | 1.75s | 5.1×     |
|                         | $10^{-8}$ | 2948    | 3.42s | 674             | 0.07s  | 674          | 5.75s | 617     | 0.13s  | <u>16</u> | 0.42s | 674                | 5.69s | <u>16</u> | 1.87s | 3.2×     |

Table 2: Comparison of iterations and runtime (in seconds) at different accuracy levels with  $\beta = 300$ . Underlined numbers indicate the iteration count of Newton methods, while non-underlined numbers correspond to the iteration count of IP-EOT/Sinkhorn.

| Case                    | Error     | IP-EOT  |        | Log-Sinkhorn |       | PSN                |       | SNS                 |       | SSNS      |       | Newton    |       | SSNS/PSN |
|-------------------------|-----------|---------|--------|--------------|-------|--------------------|-------|---------------------|-------|-----------|-------|-----------|-------|----------|
|                         |           | # iters | time   | # iters      | time  | # iters            | time  | # iters             | time  | # iters   | time  | # iters   | time  |          |
| Assignment<br>$n = 500$ | $10^{-4}$ | 26      | 0.021s | 73           | 0.26s | <u>24</u> <u>1</u> | 0.03s | <u>20</u> <u>1</u>  | 0.09s | <u>11</u> | 0.2s  | break     | 6.7×  |          |
|                         | $10^{-6}$ | 236     | 0.15s  | 654          | 2.23s | <u>24</u> <u>2</u> | 0.04s | <u>20</u> <u>5</u>  | 0.28s | <u>13</u> | 0.26s |           | 6.5×  |          |
|                         | $10^{-8}$ | 2335    | 2.93s  | 9679         | 33.6s | <u>24</u> <u>3</u> | 0.06s | <u>20</u> <u>6</u>  | 0.38s | <u>14</u> | 0.30s |           | 5×    |          |
| MNIST L2<br>$n = 784$   | $10^{-4}$ | 77      | 0.09s  | 770          | 4.68s | <u>52</u> <u>1</u> | 0.14s | <u>20</u> <u>13</u> | 1.0s  | <u>19</u> | 0.42s | <u>14</u> | 1.64s | 3.0×     |
|                         | $10^{-6}$ | 357     | 0.41s  | 1713         | 10.0s | <u>52</u> <u>3</u> | 0.20s | <u>20</u> <u>16</u> | 1.30s | <u>26</u> | 0.78s | <u>16</u> | 1.92s | 3.9×     |
|                         | $10^{-8}$ | 2776    | 4.90s  | 3128         | 18.0s | <u>52</u> <u>5</u> | 0.26s | <u>20</u> <u>17</u> | 1.41s | <u>28</u> | 0.88s | <u>17</u> | 2.05s | 7.9×     |
| MNIST L1<br>$n = 784$   | $10^{-4}$ | 111     | 0.14s  | 589          | 3.39s | <u>52</u> <u>1</u> | 0.13s | 589                 | 3.32s | 24        | 0.44s | <u>15</u> | 1.86s | 3.4×     |
|                         | $10^{-6}$ | 427     | 0.54s  | 1540         | 8.72s | <u>52</u> <u>3</u> | 0.18s | <u>700</u> <u>1</u> | 3.94s | <u>29</u> | 0.68s | <u>17</u> | 2.11s | 3.8×     |
|                         | $10^{-8}$ | 3152    | 8.94s  | 3052         | 17.3s | <u>52</u> <u>5</u> | 0.24s | <u>700</u> <u>3</u> | 4.18s | <u>31</u> | 0.78s | <u>18</u> | 2.23s | 3.3×     |

Table 3: Comparison of iterations and runtime (in seconds) at different accuracy levels and  $\beta = 200 \ln n$ . Underlined numbers indicate the iteration count of Newton methods, while non-underlined numbers correspond to the iteration count of IP-EOT/Sinkhorn. The last column reports the SSNS/PSN ratio on time, quantifying the performance gain of PSN over the second-best method.

| Case                     | Error     | IP-EOT  |       | Log-Sinkhorn |       | PSN                |       | SSNS      |       | SNS                 |       | Newton    |       | SSNS/PSN | SNS/PSN |
|--------------------------|-----------|---------|-------|--------------|-------|--------------------|-------|-----------|-------|---------------------|-------|-----------|-------|----------|---------|
|                          |           | # iters | time  | # iters      | time  | # iters            | time  | # iters   | time  | # iters             | time  | # iters   | time  |          |         |
| Assignment<br>$n = 4000$ | $10^{-4}$ | 13      | 0.69s | 8            | 2.10s | 34                 | 2.29s | <u>7</u>  | 40.2s | 8                   | 2.32s | <u>10</u> | 26.4s | 16×      | 1.0×    |
|                          | $10^{-6}$ | 120     | 6.23s | 16           | 4.08s | <u>34</u> <u>1</u> | 3.26s | <u>9</u>  | 89.8s | 16                  | 4.65s | <u>13</u> | 36.6s | 27×      | 1.4×    |
|                          | $10^{-8}$ | 1189    | 63.1s | 27           | 6.44s | <u>34</u> <u>2</u> | 4.24s | <u>12</u> | 204s  | <u>20</u> <u>1</u>  | 5.33s | <u>15</u> | 49.1s | 48×      | 1.3×    |
| MNIST L2<br>$n = 4096$   | $10^{-4}$ | 85      | 4.61s | 956          | 246s  | 68                 | 5.15s | <u>29</u> | 17.7s | <u>20</u> <u>6</u>  | 15.5s | <u>15</u> | 238s  | 3.4×     | 3.0×    |
|                          | $10^{-6}$ | 400     | 21.6s | 2133         | 489s  | <u>68</u> <u>3</u> | 7.06s | <u>42</u> | 29.4s | <u>20</u> <u>9</u>  | 23.2s | <u>19</u> | 310s  | 4.2×     | 3.3×    |
|                          | $10^{-8}$ | 693     | 37.7s | 3311         | 761s  | <u>68</u> <u>4</u> | 8.01s | <u>46</u> | 33.4s | <u>20</u> <u>31</u> | 77.6s | <u>20</u> | 328s  | 4.2×     | 9.7×    |
| MNIST L1<br>$n = 4096$   | $10^{-4}$ | 124     | 6.94s | 757          | 189s  | <u>68</u> <u>2</u> | 5.71s | <u>32</u> | 20.2s | 700                 | 164s  | <u>16</u> | 259s  | 3.5×     | 29×     |
|                          | $10^{-6}$ | 484     | 26.8s | 1974         | 444s  | <u>68</u> <u>4</u> | 7.53s | <u>42</u> | 28.0s | <u>700</u> <u>1</u> | 165s  | <u>20</u> | 328s  | 3.7×     | 22×     |
|                          | $10^{-8}$ | 3189    | 207s  | 3191         | 718s  | <u>68</u> <u>5</u> | 8.37s | <u>48</u> | 33.4s | <u>700</u> <u>4</u> | 178s  | <u>21</u> | 345s  | 4.0×     | 21×     |

Table 4: Comparison of iterations and runtime (in seconds) at different accuracy levels and  $\beta = 200 \ln n$ . Underlined numbers indicate the iteration count of Newton methods, while non-underlined numbers correspond to the iteration count of IP-EOT/Sinkhorn.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

## J DISCUSSION OF THE SWITCHING PARAMETER

The choice of  $\lambda$  aims to select the more efficient option between Proximal-Sinkhorn (PS) and Proximal-Newton (PN). Accordingly, this section presents the efficiency of PS and PN under different values of  $\beta$ , together with the corresponding changes in the sparsity of  $\hat{H}$ , providing guidance for choosing an appropriate  $\lambda$ .

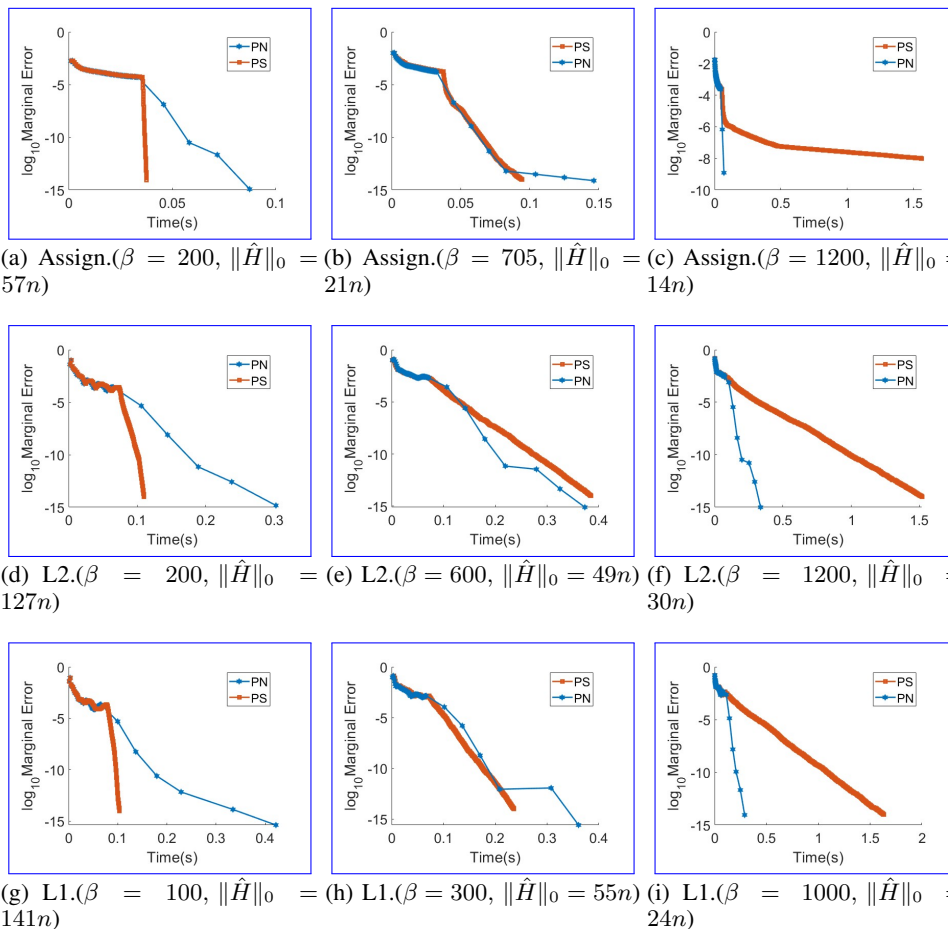


Figure 9: Performance of PS and PN algorithm on three datasets with different  $\beta$  value.

To discuss the selection of  $\lambda$ , we evaluate PS and PN under different values of  $\beta$  to explore the relation between the number of non-zero elements  $|\hat{H}|_0$  and performance, we specifically selected three sets of  $\beta$  such that the three scenarios—namely, PS being superior, the two being comparable, and PN being superior—are clearly demonstrated. As shown in Fig. 9, as  $\beta$  increases, the number of non-zero elements  $|\hat{H}|_0$  gradually decreases, and PN progressively outperforms PS. Interestingly, under both L1 and L2 settings, when PS and PN achieve comparable time performance, their sparsity levels are similar even though the values of  $\beta$  differ, which validates the effectiveness of using  $\lambda$  as the switching parameter.

## K COMPARISON WITH EXISTING FIRST ORDER METHODS

To illustrate the efficiency of our Proximal-Sinkhorn method (PS), we compare it with two classical entropic solvers: log-Sinkhorn and Epsilon-Scaling method<sup>5</sup>. We compare the algorithms for  $\beta$  in the range (100, 700), since for larger values second-order methods gain a clear advantage (as shown in Figure 9).

As shown in Figure 10, the error curve of the epsilon-scaling method exhibits an initial plateau before decreasing, because it requires a warm-up phase. Our PS method does not suffer from this issue and consistently achieves superior performance.

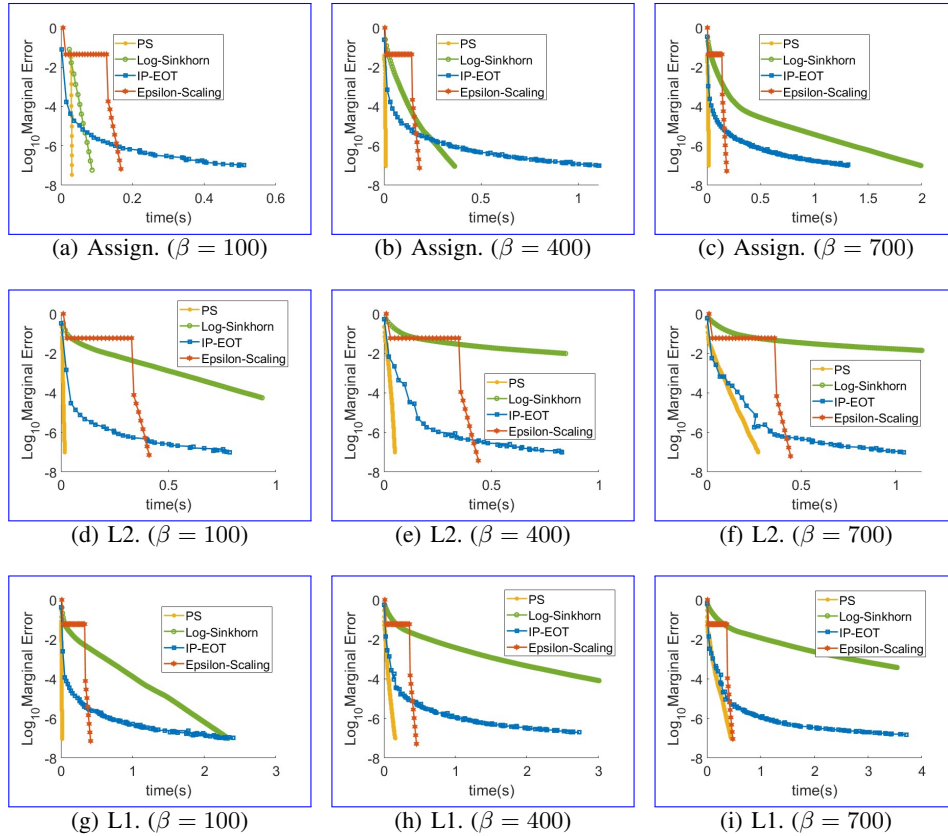


Figure 10: Comparison of four first-order algorithms.

## L EVALUATION ON IMAGENET

We consider an OT problem between two categories of images from the ImageNet dataset. We use a subset of ImageNet from the Imagenette GitHub repository<sup>6</sup>, which contains ten image classes. Approximately 1000 images per class are selected. Each image is mapped to a 30-dimensional feature vector by first passing it through a ResNet18 network to obtain a 512-dimensional embedding, followed by dimensionality reduction via principal component analysis.

Let  $x_i \in \mathbb{R}^{30}$  denote the feature vector of an image in the first category for  $i = 1, \dots, n$ , and let  $y_j \in \mathbb{R}^{30}$  denote the feature vector of an image in the second category for  $j = 1, \dots, m$ . We set  $\mathbf{a} = n^{-1}\mathbf{1}_n$ ,  $\mathbf{b} = m^{-1}\mathbf{1}_m$ , and define the cost matrix either as  $C_{ij} = \|x_i - y_j\|_1$  or  $C_{ij} = \|x_i - y_j\|_2$ . The problem size is  $n \approx 1000$ .

<sup>5</sup>Schmitzer, B. (2016). Stabilized Sparse Scaling Algorithms for Entropy Regularized Transport Problems. arXiv:1610.06519.

<sup>6</sup><https://github.com/fastai/imagenette>

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

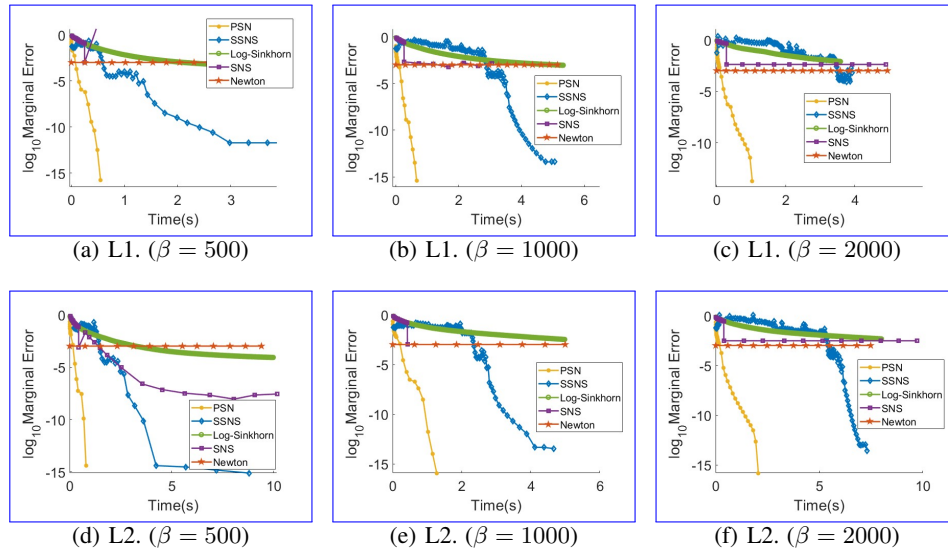


Figure 11: Performance on the ImageNet dataset.

As shown in Figure 11, both Newton and SNS encounter failures in at least five cases, and even SSNS fails in one instance. In contrast, PSN achieves a perfect success rate across all test cases, demonstrating reliable performance and high efficiency in every setting.

## M DISCUSSION OF THE PROXIMAL STEP SIZE

This section examines the sensitivity of PSN to the proximal step size  $\Delta\beta$ . According to Theorem 1, the proximal step size  $\Delta\beta$  influences only the initialization of the second stage. When the second stage is achieved by Sinkhorn scaling,  $\Delta\beta$  does not affect the convergence rate. Therefore, this focuses on the sensitivity of the Proximal-Newton (PN) method to the choice of  $\Delta\beta$ .

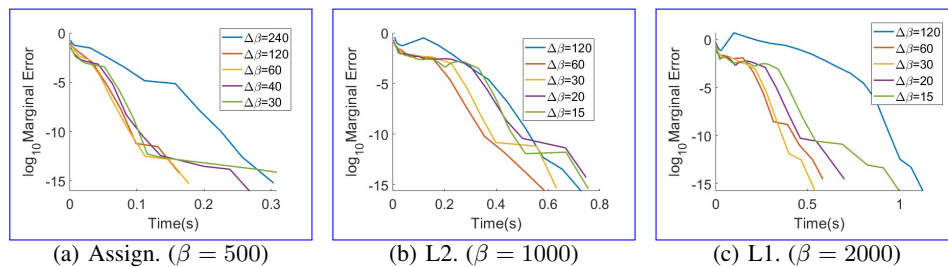
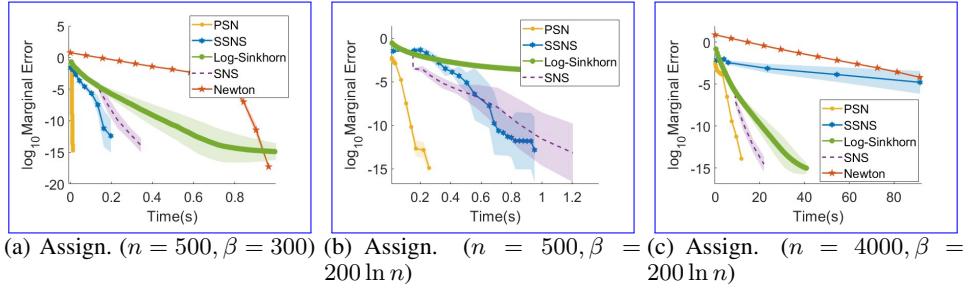


Figure 12: Impact of the proximal step size  $\Delta\beta$  on proximal-newton method.

As illustrated in Fig. 12, the PN algorithm exhibits considerable robustness to the choice of proximal step size, consistently achieving stable convergence across different datasets and values of  $\beta$ . The convergence behavior remains highly similar across a wide range of step sizes, except when  $\Delta\beta$  is excessively large—that is, when the number of IP-EOT iterations  $\ell$  is particularly small. Even in such cases, the algorithm only suffers from a slightly slower initial convergence rate, with the total computation time remaining comparable. On the other hand, when  $\beta$  is large and the proximal step size is small (corresponding to a large number of IP-EOT iterations  $\ell$ ), the computational cost of IP-EOT increases somewhat. Overall, as long as the proximal step size is chosen to ensure the numerical stability of the IP-EOT phase, the algorithm PSN achieves favorable convergence performance.

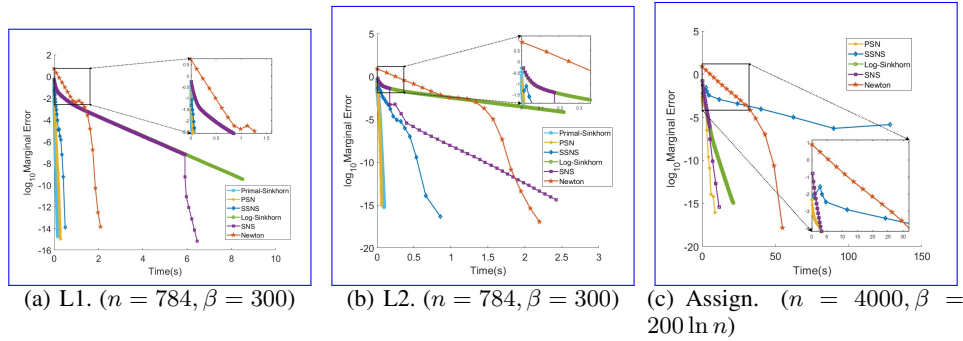
1458 **N** VARIABILITY OF THE METHODS UNDER DIFFERENT RANDOM  
 1459 REALIZATIONS  
 1460

1461 For each scenario presented in the Figure 5 (a), (d) and (g), we conducted twenty independent  
 1462 random trials and plotted the results with variance. Figure 13 demonstrates that Newton and PSN  
 1463 exhibit small fluctuations under randomness, whereas SNS and SSNS show much larger variability.  
 1464



1476 Figure 13: Performance under random realizations on the linear assignment tasks.  
 1477

1480 **O** ZOOM-IN OF SOME FIGURES  
 1481



1495 Figure 14: Zoom-in plots.  
 1496

1497  
1498 **P** TIME PER ITERATION  
 1499

| Case                    | IP-EOT  | Log-Sinkhorn | PS      | Epsilon-Scaling |
|-------------------------|---------|--------------|---------|-----------------|
| Assignment<br>$n = 500$ | $5e-4s$ | $5e-3s$      | $6e-4s$ | $4e-3s$         |
| MNIST L2<br>$n = 784$   | $1e-3s$ | $8e-3s$      | $6e-4s$ | $0.014s$        |
| MNIST L1<br>$n = 784$   | $1e-3s$ | $9e-3s$      | $6e-4s$ | $0.014s$        |

1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507 Table 5: Time per iteration of four first-order methods under  $\beta = 300$ .  
 1508  
1509

1510  
1511 **Q** SPECTRAL ESTIMATION

| Case                    | PN (Newton part) | SSNS  | SNS (Newton part) | Newton |
|-------------------------|------------------|-------|-------------------|--------|
| Assignment<br>$n = 500$ | 0.03s            | 0.04s | 0.1s              | break  |
| MNIST L2<br>$n = 784$   | 0.04s            | 0.03s | 0.2s              | 0.1s   |
| MNIST L1<br>$n = 784$   | 0.03s            | 0.03s | 0.3s              | 0.15s  |

Table 6: Time per iteration of four second-order methods under  $\beta = 1200$ .

By invoking Gershgorin’s Circle Theorem, the eigenvalues of the Hessian matrix  $H$ , defined in equation (15), satisfy

$$\lambda(H) \in \left( \bigcup_{i=1}^n [0, 2\beta \mathbf{a}_i] \right) \cup \left( \bigcup_{j=1}^n [0, 2\beta \mathbf{b}_j] \right) = [0, 2\beta M_{max}], \quad (109)$$

where  $M = \max_{i,j}(\mathbf{a}_i, \mathbf{b}_j)$ . Let  $\Delta a_i = a_i - \sum_j \hat{P}_{ij}$  be the source mass loss, analogously for  $\Delta b_j = b_j - \sum_i \hat{P}_{ij}$  as target mass loss. Let  $\Delta_{\min} = \min_{i,j}(\Delta a_i, \Delta b_j)$ , the eigenvalues  $\lambda$  of  $\hat{H}$  satisfy the following bounds

$$\beta \Delta_{\min} \leq \lambda(\hat{H}) \leq 2\beta M - \beta \Delta_{\min}. \quad (110)$$

To illustrate the behavior under different problem scales, we use the linear assignment problem to validate the range.

**Upper bound.** In the context of the linear assignment problem, we set  $\mathbf{a} = \mathbf{b} = \mathbf{1}/n$ , which yields an upper bound of  $2\beta/n$  for the largest eigenvalue. As illustrated in Figure 15, the largest eigenvalue of  $\hat{H}$  consistently adheres to this theoretical bound across all tested scenarios.

**Lower bound.** As illustrated in Figure 16, the smallest eigenvalue reaches a steady state within a negligible number of initial iterations and remains virtually constant thereafter, regardless of the problem size  $n$  or the regularization parameter  $\beta$ . It demonstrates that our method effectively bounds the smallest eigenvalue away from zero, thereby ensuring numerical stability across a broad parameter space.

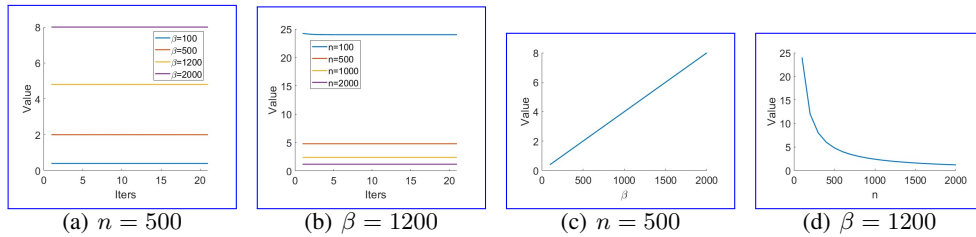


Figure 15: The largest eigenvalues of  $\hat{H}$  with different  $n$  and  $\beta$ .

As for the feature dimension  $d$ , Figure 17 show that the upper bound of the eigenvalues is almost unaffected by  $d$ , while the lower bound experiences slight perturbations.

Moreover, both the largest and smallest eigenvalues of  $\hat{H}$  are almost independent of the iteration, which indirectly confirms the robustness of our sparse Newton method.

## R LARGER SCALE PROBLEM

We compared our algorithm with other algorithms on the linear assignment dataset with  $n = 10^4$ . Figure 18 demonstrates the superiority of our PSN.

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

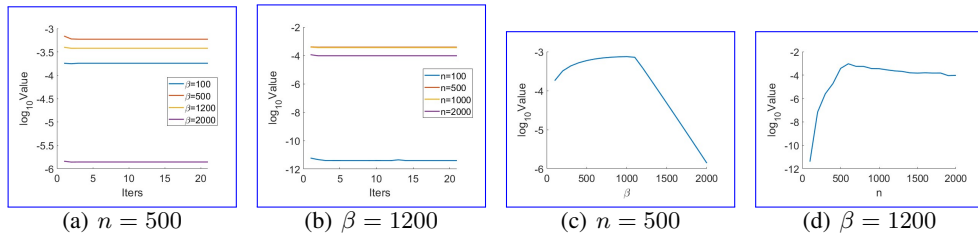


Figure 16: The smallest eigenvalues of  $\hat{H}$  with different  $n$  and  $\beta$ .

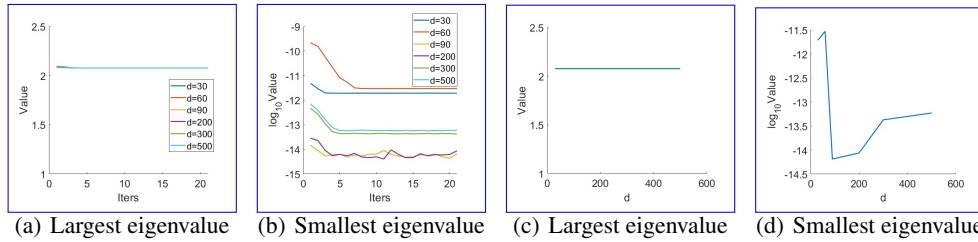


Figure 17: The largest and smallest eigenvalue of  $\hat{H}$  with different feature dimension  $d$  under  $n = 963$  and  $\beta = 1000$ .

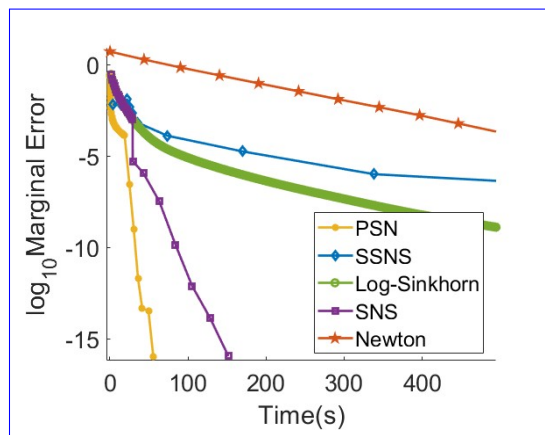


Figure 18: Assign. ( $n = 10000$ )

---

1620 THE USE OF LARGE LANGUAGE MODEL (LLM)  
1621

1622 We used a large language model (LLM) solely for grammar checking and stylistic polishing of the  
1623 manuscript text. The LLM did not generate scientific content, ideas, experiments, analyses, code,  
1624 or results. All technical contributions, claims, and conclusions are the authors' own. The authors  
1625 verified all edited text and remain fully responsible for the content. No confidential, proprietary, or  
1626 personally identifying information was provided to the LLM.  
1627

1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673