
GROMOV-WASSERSTEIN LEARNING FOR GRAPH MATCHING AND NODE EMBEDDING

Hongteng Xu^{1,2} Dixin Luo² Hongyuan Zha³ Lawrence Carin²

OVERVIEW

- **Introduction**
- **Gromov-Wasserstein Learning (GWL) framework**
- **Experiments**

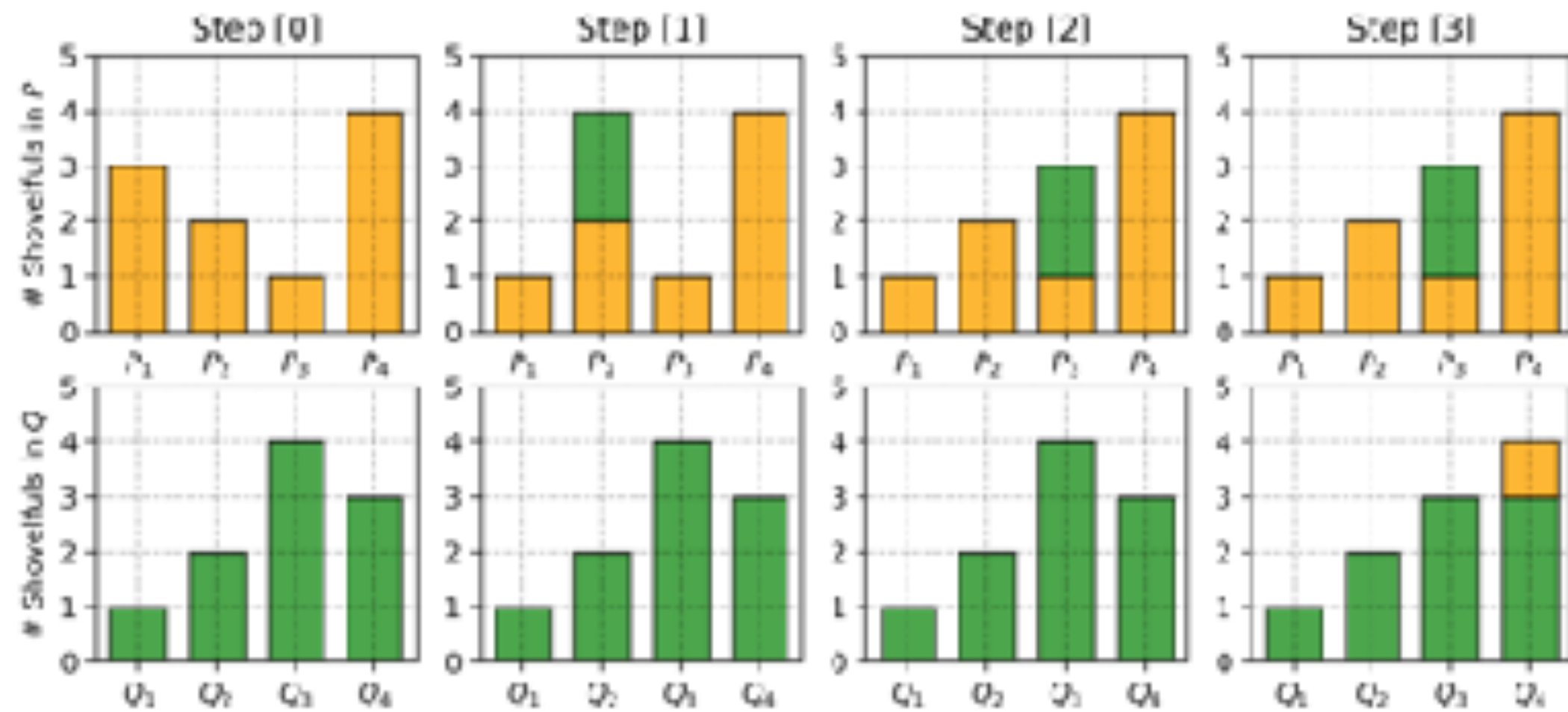


WASSERSTEIN DISTANCE

- **K-Wasserstein distance:** $(WD - k) \quad \inf \left\{ K(\gamma) := \int_{X \times Y} \|x - y\|_k^k d\gamma(x, y) : \gamma \in \Pi(\mu, \nu) \right\}$

- **Earth mover (1-wasserstein distance):** $(WD - 1) \quad \inf \left\{ K(\gamma) := \int_{X \times Y} \|x - y\|_1 d\gamma(x, y) : \gamma \in \Pi(\mu, \nu) \right\}$

➤ $W(\mathbb{P}_p, \mathbb{P}_q) = \inf_{\gamma \in \Pi(\mathbb{P}_p, \mathbb{P}_q)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|] = \min_{\mathbf{T} \in \Pi(\mathbf{p}, \mathbf{q})} \sum_{i=1}^n \sum_{j=1}^m \mathbf{T}_{ij} \cdot c(\mathbf{x}_i, \mathbf{y}_j)$



$w = 2 * d(P1, Q2) + 2 * d(P2, Q2) + 1 * d(P3, Q4) = 5$

P, Q	Q1	Q2	Q3	Q4
P1	0	2	0	0
P2	0	2	0	0
P3	0	0	0	1
P4	0	0	0	0

GROMOV-WASSERSTEIN DISTANCE

Definition 2.1. Let (X, d_X, μ_X) and (Y, d_Y, μ_Y) be two metric measure spaces, where (X, d_X) is a compact metric space and μ_X is a Borel probability measure on X (with (Y, d_Y, μ_Y) defined in the same way). The Gromov-Wasserstein distance $d_{GW}(\mu_X, \mu_Y)$ is defined as:

Better explain the terms one by one and slowly.

Sometimes you need to read some tutorials to understand the terms.

$$\inf_{\pi \in \Pi(\mu_X, \mu_Y)} \iint_{X \times Y, X \times Y} L(x, y, x', y') d\pi(x, y) d\pi(x', y')$$

where $L(x, y, x', y') = |d_X(x, x') - d_Y(y, y')|$ is the loss function and $\Pi(\mu_X, \mu_Y)$ is the set of all probability measures on $X \times Y$ with μ_X and μ_Y as marginals.

PROPOSED METHOD

- ▶ **Graph matching:** Finding a correspondence between their nodes.
- ▶ **Node embedding:** Embedding their nodes in the same space.

Unify them in the **Gromov-Wasserstein Learning (GWL)** framework.

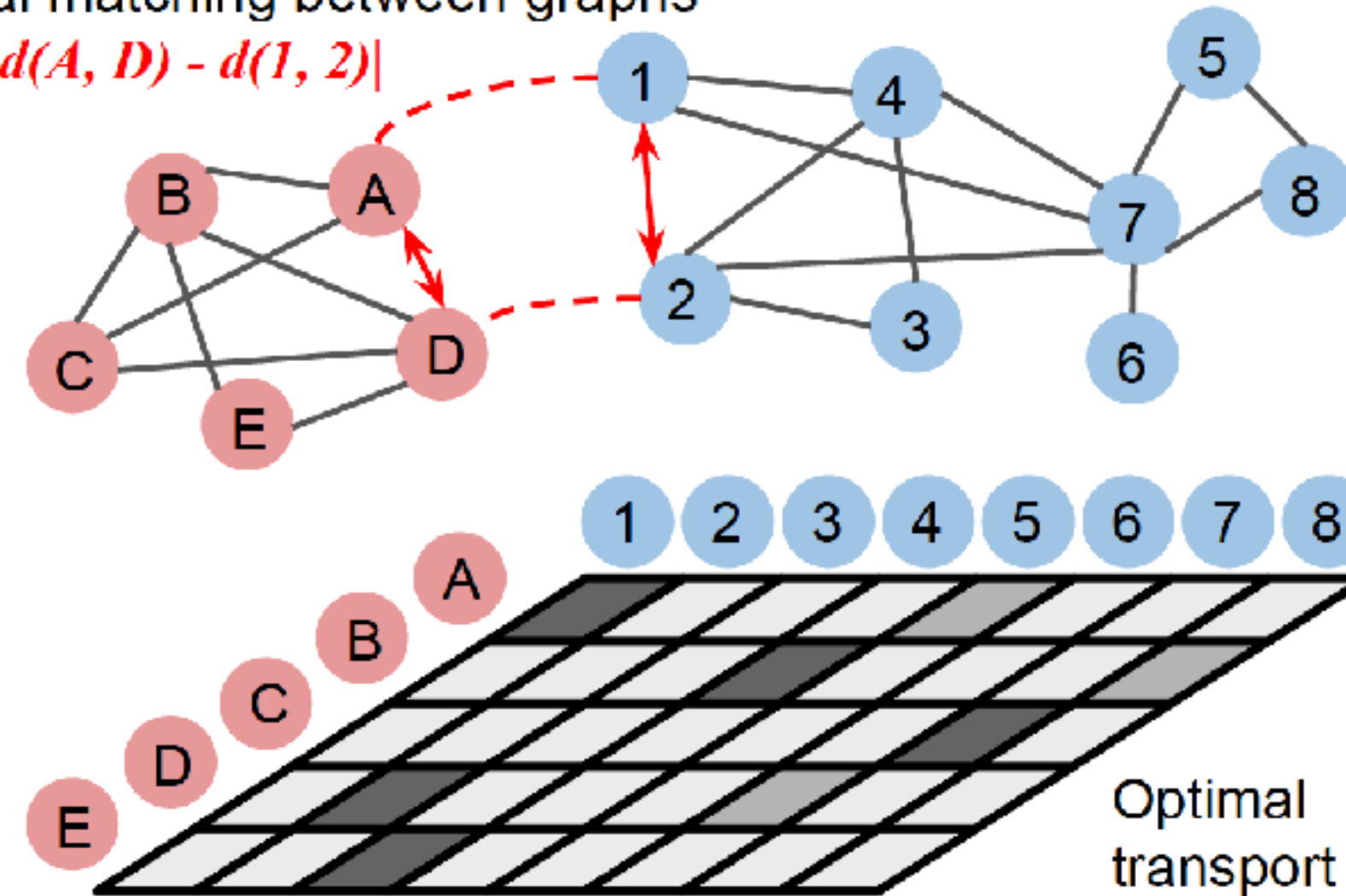
better explain the terms using the example.

$$d_{GW}(G_s, G_t) := \min_{\mathbf{T} \in \Pi(\mu_s, \mu_t)} \sum_{i,j,i',j'} L(c_{ij}^s, c_{i'j'}^t) T_{ii'} T_{jj'} = \min_{\mathbf{T} \in \Pi(\mu_s, \mu_t)} \langle \mathbf{L}(\mathbf{C}_s, \mathbf{C}_t, \mathbf{T}), \mathbf{T} \rangle.$$

$$\Pi(\mu_s, \mu_t) = \{ \mathbf{T} \in \mathbb{R}^{|\mathcal{V}_s| \times |\mathcal{V}_t|} \mid \mathbf{T} \mathbf{1}_{|\mathcal{V}_t|} = \mu_s, \mathbf{T}^\top \mathbf{1}_{|\mathcal{V}_s|} = \mu_t \}$$

Relational matching between graphs

$$\text{Cost} = |d(A, D) - d(1, 2)|$$



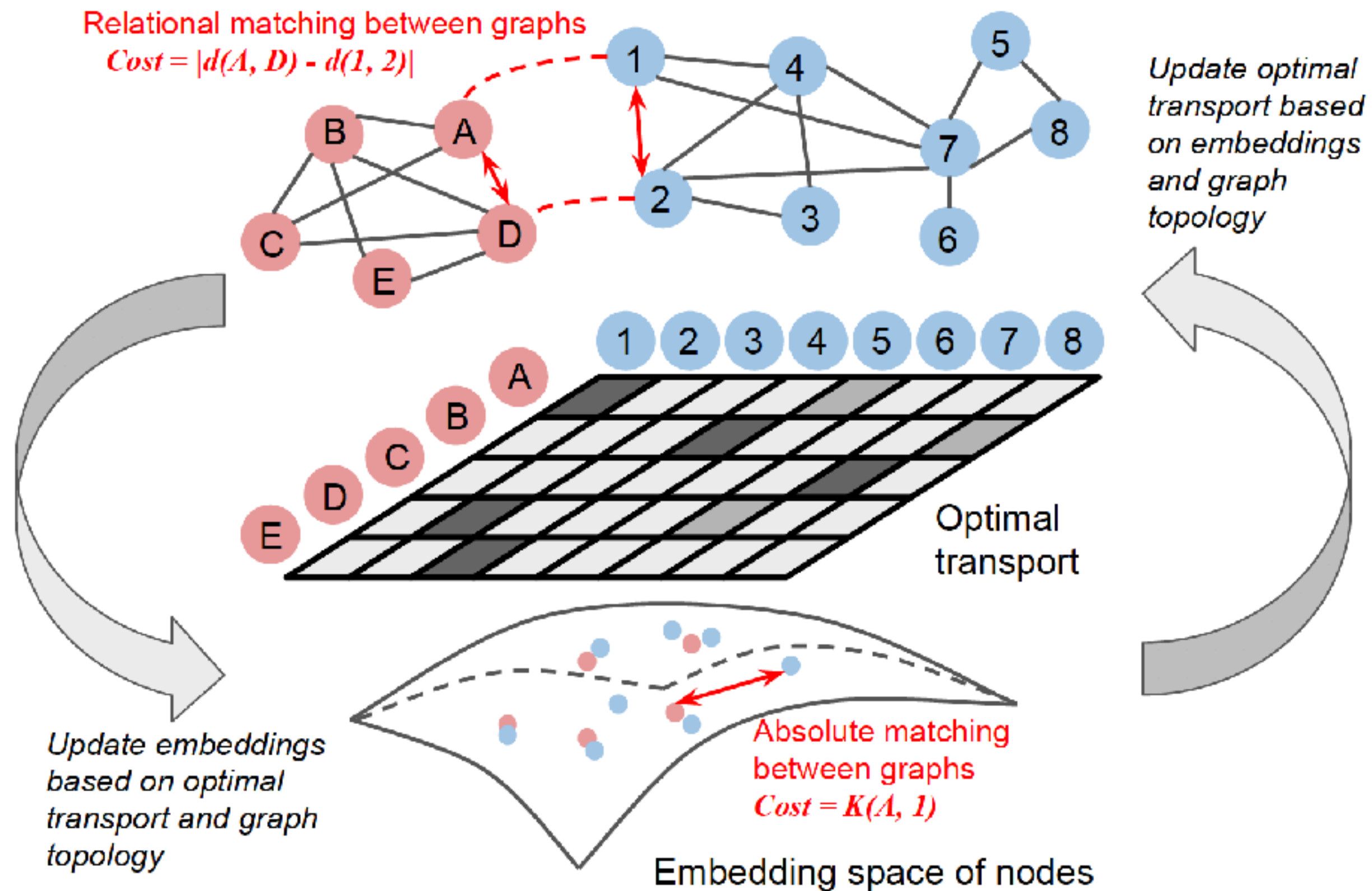
$$L_{jj'} = \sum_{i,i'} L(c_{ij}^s, c_{i'j'}^t) T_{ii'}$$

$$\mathbf{L}(\mathbf{C}_s, \mathbf{C}_t, \mathbf{T}) = [L_{jj'}] \in \mathbb{R}^{|\mathcal{V}_s| \times |\mathcal{V}_t|}$$

there is a lot of information and math terms.
Better organize them by input-output and objective.

GROMOV-WASSERSTEIN LEARNING

$$\min_{\mathbf{X}_s, \mathbf{X}_t} \min_{\mathcal{T} \in \Pi(\mu_s, \mu_t)} \underbrace{\langle \mathbf{L}(\mathbf{C}_s(\mathbf{X}_s), \mathbf{C}_t(\mathbf{X}_t), \mathcal{T}), \mathcal{T} \rangle}_{\text{Gromov-Wasserstein discrepancy}} + \underbrace{\alpha \langle \mathbf{K}(\mathbf{X}_s, \mathbf{X}_t), \mathcal{T} \rangle}_{\text{Wasserstein discrepancy}} + \underbrace{\beta \sum_{k=s,t} R(\mathbf{K}(\mathbf{X}_k, \mathbf{X}_k), \mathbf{C}_k)}_{\text{prior information}}.$$



GROMOV-WASSERSTEIN LEARNING

$$\min_{\mathbf{X}_s, \mathbf{X}_t} \min_{\mathbf{T} \in \Pi(\mu_s, \mu_t)} \underbrace{\langle \mathbf{L}(\mathbf{C}_s(\mathbf{X}_s), \mathbf{C}_t(\mathbf{X}_t), \mathbf{T}), \mathbf{T} \rangle}_{\text{Gromov-Wasserstein discrepancy}} + \underbrace{\alpha \langle \mathbf{K}(\mathbf{X}_s, \mathbf{X}_t), \mathbf{T} \rangle}_{\text{Wasserstein discrepancy}} + \underbrace{\beta \sum_{k=s,t} R(\mathbf{K}(\mathbf{X}_k, \mathbf{X}_k), \mathbf{C}_k)}_{\text{prior information}}.$$

- $\mathbf{C}_k(\mathbf{X}_k) = (1 - \alpha)\mathbf{C}_k + \alpha\mathbf{K}(\mathbf{X}_k, \mathbf{X}_k)$, for $k = s, t$
 - $\mathbf{K}(\mathbf{X}_k, \mathbf{X}_k) = [\kappa(\mathbf{x}_i^k, \mathbf{x}_j^k)] \in \mathbb{R}^{|\mathcal{V}_k| \times |\mathcal{V}_k|}$ **distance between the node embedding within the same graph**
 - $\mathbf{K}(\mathbf{X}_s, \mathbf{X}_t) = [\kappa(\mathbf{x}_i^s, \mathbf{x}_j^t)] \in \mathbb{R}^{|\mathcal{V}_s| \times |\mathcal{V}_t|}$ **distance between nodes of the two graphs**
 - $R(\mathbf{X}_s, \mathbf{X}_t) = \sum_{k=s,t} L(\mathbf{K}(\mathbf{X}_k, \mathbf{X}_k), \mathbf{C}_k) + \underbrace{L(\mathbf{K}(\mathbf{X}_s, \mathbf{X}_t), \mathbf{C}_{st})}_{\text{optional}}$
-

GROMOV-WASSERSTEIN LEARNING

➤ **Updating optimal transport: m-th outer iteration, n-th inner iteration, based on KL divergence**

➤ $\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t)} \langle \mathbf{L}(\mathbf{C}_s(\mathbf{X}_s^{(m)}), \mathbf{C}_t(\mathbf{X}_t^{(m)}), \mathbf{T}), \mathbf{T} \rangle + \alpha \langle \mathbf{K}(\mathbf{X}_s^{(m)}, \mathbf{X}_t^{(m)}), \mathbf{T} \rangle + \gamma \text{KL}(\mathbf{T} \parallel \mathbf{T}^{(n)}). \quad (6)$

➤ $\text{KL}(\mathbf{T} \parallel \mathbf{T}^{(n)}) = \sum_{ij} T_{ij} \log \frac{T_{ij}}{T_{ij}^{(n)}} - T_{ij} + T_{ij}^{(n)}$

➤ **Parameter α control:** $\alpha_m = \frac{m}{M}$

➤ **Updating embeddings:** $\min_{\mathbf{X}_s, \mathbf{X}_t} \alpha_m \langle \mathbf{K}(\mathbf{X}_s, \mathbf{X}_t), \hat{\mathbf{T}}^{(m)} \rangle + \beta R(\mathbf{X}_s, \mathbf{X}_t). \quad (8)$



EMBEDDING-BASED DISTANCE

➤ **Cosine-based distance:** $k(\mathbf{x}_i, \mathbf{x}_j) = 1 - \exp\left(-\sigma\left(1 - \frac{\mathbf{x}_i^\top \mathbf{x}_j}{\|\mathbf{x}_i\|_2 \|\mathbf{x}_j\|_2}\right)\right)$

➤ **Radial basis function(RBF)-based distance:** $k(\mathbf{x}_i, \mathbf{x}_j) = 1 - \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}\right)$



ALGORITHM

$$\min_{T \in \Pi(\mu_s, \mu_t)} \langle L(C_s(\mathbf{X}_s^{(m)}), C_t(\mathbf{X}_t^{(m)}), T), T \rangle + \alpha \langle K(\mathbf{X}_s^{(m)}, \mathbf{X}_t^{(m)}), T \rangle + \gamma \text{KL}(T \| T^{(n)}). \quad (6)$$

$$\min_{\mathbf{X}_s, \mathbf{X}_t} \alpha_m \langle K(\mathbf{X}_s, \mathbf{X}_t), \hat{T}^{(m)} \rangle + \beta R(\mathbf{X}_s, \mathbf{X}_t). \quad (8)$$

Algorithm 1 Gromov-Wasserstein Learning (GWL)

- 1: **Input:** $\{C_s, C_t\}$, $\{\mu_s, \mu_t\}$, β, γ , the dimension D , the number of outer/inner iterations $\{M, N\}$.
 - 2: **Output:** $\mathbf{X}_s, \mathbf{X}_t$ and \hat{T} .
 - 3: Initialize $\mathbf{X}_s^{(0)}, \mathbf{X}_t^{(0)}$ randomly, $\hat{T}^{(0)} = \mu_s \mu_t^\top$.
 - 4: **For** $m = 0 : M - 1$
 - 5: Set $\alpha_m = \frac{m}{M}$.
 - 6: **For** $n = 0 : N - 1$
 - 7: Update optimal transport $\hat{T}^{(m+1)}$ via solving (6).
 - 8: Obtain $\mathbf{X}_s^{(m+1)}, \mathbf{X}_t^{(m+1)}$ via solving (8).
 - 9: $\mathbf{X}_s = \mathbf{X}_s^{(M)}, \mathbf{X}_t = \mathbf{X}_t^{(M)}$ and $\hat{T} = \hat{T}^{(M)}$.
 - 10: $\backslash\backslash$ Graph matching:
 - 11: Initialize correspondence set $\mathcal{P} = \emptyset$
 - 12: **For** $v_i \in \mathcal{V}_s$
 - 13: $j = \arg \max_j \hat{T}_{ij}$. $\mathcal{P} = \mathcal{P} \cup \{(v_i \in \mathcal{V}_s, v_j \in \mathcal{V}_t)\}$.
-

SINKHORN-KNOPP ALGORITHM

➤ $\mathcal{L}(P, \alpha, \beta) = \sum_{ij} \gamma p_{ij} \log p_{ij} + p_{ij} c_{ij} + \alpha^T (P \mathbf{1}_d - \mathbf{a}) + \beta^T (P^T \mathbf{1}_d - \mathbf{b})$

➤ $\frac{\partial \mathcal{L}}{\partial p_{ij}} = 0$, then $p_{ij} = e^{-\frac{\alpha_i}{\gamma}} e^{-\frac{c_{ij}}{\gamma}} e^{-\frac{\beta_j}{\gamma}}$

➤ $P = \text{diag}(u) K \text{diag}(v) = \text{diag}(u) e^{-\frac{c_{ij}}{\gamma}} \text{diag}(v)$

➤ $\mathbf{u} \odot (\mathbf{K} \mathbf{v}) = \mathbf{a}$ and $\mathbf{v} \odot (\mathbf{K}^T \mathbf{u}) = \mathbf{b}$

➤ $\mathbf{u}^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{a}}{\mathbf{K} \mathbf{v}^{(\ell)}}$ and $\mathbf{v}^{(t+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{b}}{\mathbf{K}^T \mathbf{u}^{(\ell+1)}}$

PROXIMAL POINT METHOD

$$\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}_s, \boldsymbol{\mu}_t)} \langle \mathbf{C}^{(m,n)} - \gamma \log \mathbf{T}^{(n)}, \mathbf{T} \rangle + \gamma H(\mathbf{T}), \quad (7)$$

$$\text{where } \mathbf{C}^{(m,n)} = \mathbf{L}(\mathbf{C}_s, \mathbf{C}_t, \mathbf{T}^{(n)}) + \alpha \mathbf{K}(\mathbf{X}_s^{(m)}, \mathbf{X}_t^{(m)}) + \gamma,$$

$$\text{and } H(\mathbf{T}) = \sum_{i,j} T_{ij} \log T_{ij}$$

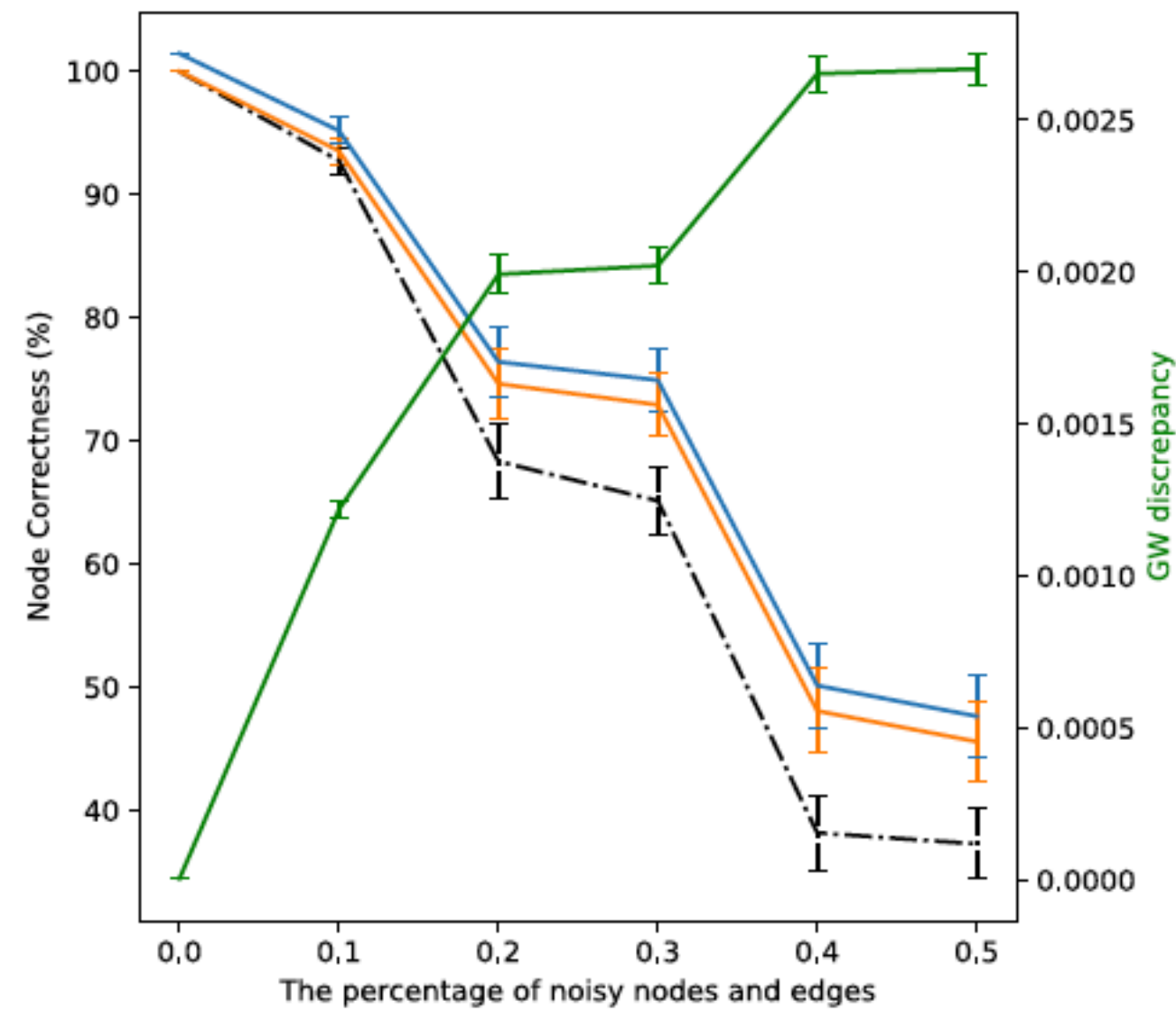
Algorithm 2 Proximal Point Method for GW Discrepancy

- 1: **Input:** $\{\mathbf{C}_s, \mathbf{C}_t\}$, $\{\boldsymbol{\mu}_s, \boldsymbol{\mu}_t\}$, current embeddings $\{\mathbf{X}_s^{(m)}, \mathbf{X}_t^{(m)}\}$, γ , the number of inner iterations N .
 - 2: **Output:** $\hat{\mathbf{T}}^{(m+1)}$.
 - 3: Initialize $\mathbf{T}^{(0)} = \boldsymbol{\mu}_s \boldsymbol{\mu}_t^\top$ and $\mathbf{a} = \boldsymbol{\mu}_s$.
 - 4: **for** $n = 0 : N - 1$ **do**
 - 5: Calculate the $\mathbf{C}^{(m,n)}$ in (7).
 - 6: Set $\mathbf{G} = \exp(-\frac{\mathbf{C}^{(m,n)}}{\gamma}) \odot \mathbf{T}^{(n)}$.
 - 7: \\ Sinkhorn-Knopp algorithm:
 - 8: **for** $j = 1 : J$ **do**
 - 9: $\mathbf{b} = \frac{\boldsymbol{\mu}_Y}{\mathbf{G}^\top \mathbf{a}}$
 - 10: $\mathbf{a} = \frac{\boldsymbol{\mu}_X}{\mathbf{G} \mathbf{b}}$
 - 11: **end for**
 - 12: $\mathbf{T}^{(n+1)} = \text{diag}(\mathbf{a}) \mathbf{G} \text{diag}(\mathbf{b})$.
 - 13: **end for**
 - 14: $\hat{\mathbf{T}}^{(m+1)} = \mathbf{T}^{(N)}$.
-

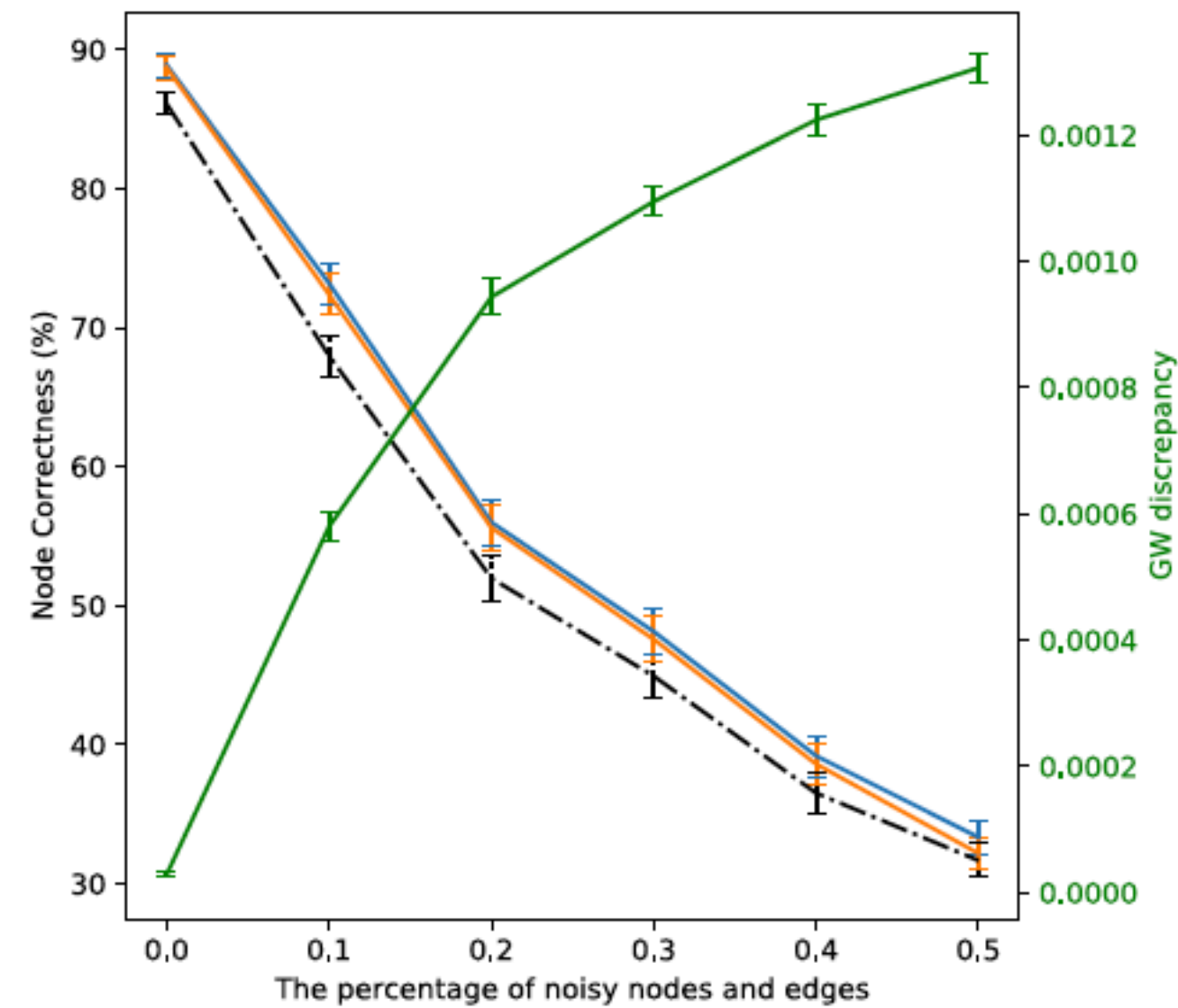
EXPERIMENTS

Explain the x and y axes first.
Then explain the curves.
Lastly, explain the settings and the performance.

--- GWD — GWL-C (OT) — GWL-C (Embedding) — GW Discrepancy



(a) K-NN: $|\mathcal{V}_s| = 50$

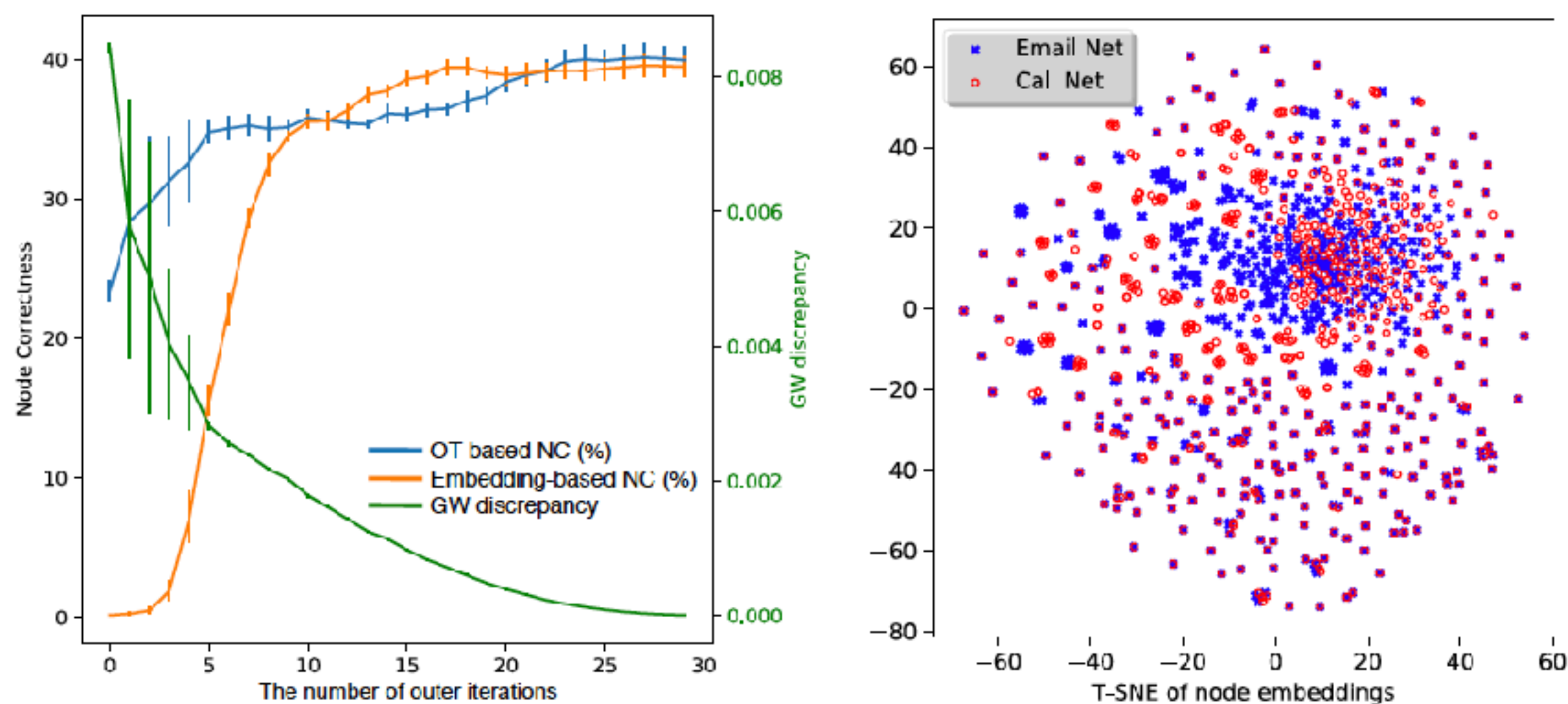


(b) BA: $|\mathcal{V}_s| = 50$

Figure 2. The performance of our method on synthetic data.

EXPERIMENTS

need to explain the baselines,
and the differences between
the baselines and the proposed method
how the difference lead to the
the performance improvement.



(a) Stability and convergence (b) Learned embeddings

Figure 3. Visualization of typical experimental results.

Table 1. Communication network matching results.

Method	Call→Email (Sparse)	Call→Email (Dense)
	Node Correctness (%)	Node Correctness (%)
GAA	34.22	0.53
LRSA	38.20	2.93
TAME	37.39	2.67
GRAAL	39.67	0.48
MI-GRAAL	35.53	0.64
MAGNA++	7.88	0.09
HugAlign	36.21	3.86
NETAL	36.87	1.77
GWD	23.16±0.46	1.77±0.22
GWL-R	39.64±0.57	3.80±0.23
GWL-C	40.45±0.53	4.23±0.27

REFERENCE

- **Cuturi, M. and Peyr e, G. Computational optimal transport.2017.**
 - **Peyr e, G., Cuturi, M., and Solomon, J. Gromov-wasserstein averaging of kernel and distance matrices. In ICML, 2016.**
 - **Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. In NIPS, pp. 2292–2300, 2013.**
 - **Benamou, J.-D., Carlier, G., Cuturi, M., Nenna, L., and Peyr e, G. Iterative Bregman projections for regularized transportation problems. SIAM Journal on Scientific Computing, 37(2):A1111–A1138, 2015.**
-