

# Fused Gromov-Wasserstein barycenter with application on graphs

Zhichen Zeng

10/13

The logo of the University of Illinois, featuring a red "I" followed by the word "ILLINOIS" in blue capital letters.A grayscale photograph of the University of Illinois dome, showing the ornate architectural details of the roof and the central cupola.

# Outline



- Motivation
- Preliminaries
- Fused Gromov-Wasserstein barycenter
- Experiments
- Takeaways

# Outline



- Motivation
- Preliminaries
- Fused Gromov-Wasserstein barycenter
- Experiments
- Takeaways

# Motivation



- Networks are everywhere
  - Alignment (node-level)
  - Clustering (subgraph-level)
  - Comparison (graph-level)
- Distance measures at different levels are important



# Motivation



- Common distance measures
  - Invalid in certain cases
  - Ignore underlying structure

$$\text{KL divergence: } KL(P||Q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

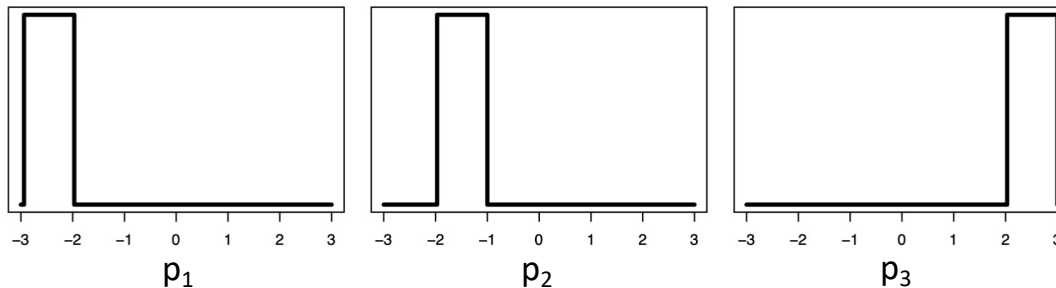
$$\text{JS divergence: } JS(P||Q) = \frac{1}{2} KL(P||\frac{P+Q}{2}) + \frac{1}{2} KL(Q||\frac{P+Q}{2})$$

$$\text{Total variation: } \frac{1}{2} \int |p(x) - q(x)| dx$$

$$\text{Hellinger distance: } \sqrt{\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx}$$

$$L_2 \text{ distance: } \int (p(x) - q(x))^2 dx$$

$$\chi_2 \text{ distance: } \int \frac{(p(x) - q(x))^2}{q(x)} dx$$



$KL(p_1, p_2)$ ,  $KL(p_1, p_3)$ : invalid

$JS(p_1, p_2)$ ,  $JS(p_1, p_3)$ : invalid

$$\frac{1}{2} \int |p_1(x) - p_2(x)| dx = \frac{1}{2} \int |p_1(x) - p_3(x)| dx = 1$$

- Optimal transport and Wasserstein distance

[1] Peyré, Gabriel, and Marco Cuturi. "Computational Optimal Transport." *arXiv preprint arXiv:1803.00567* (2018).

# Outline



- Motivation
- Preliminaries
  - Optimal transport and Wasserstein distance
  - Gromov-Wasserstein distance
  - Fused Gromov-Wasserstein distance
- Fused Gromov-Wasserstein barycenter
- Experiments
- Takeaways

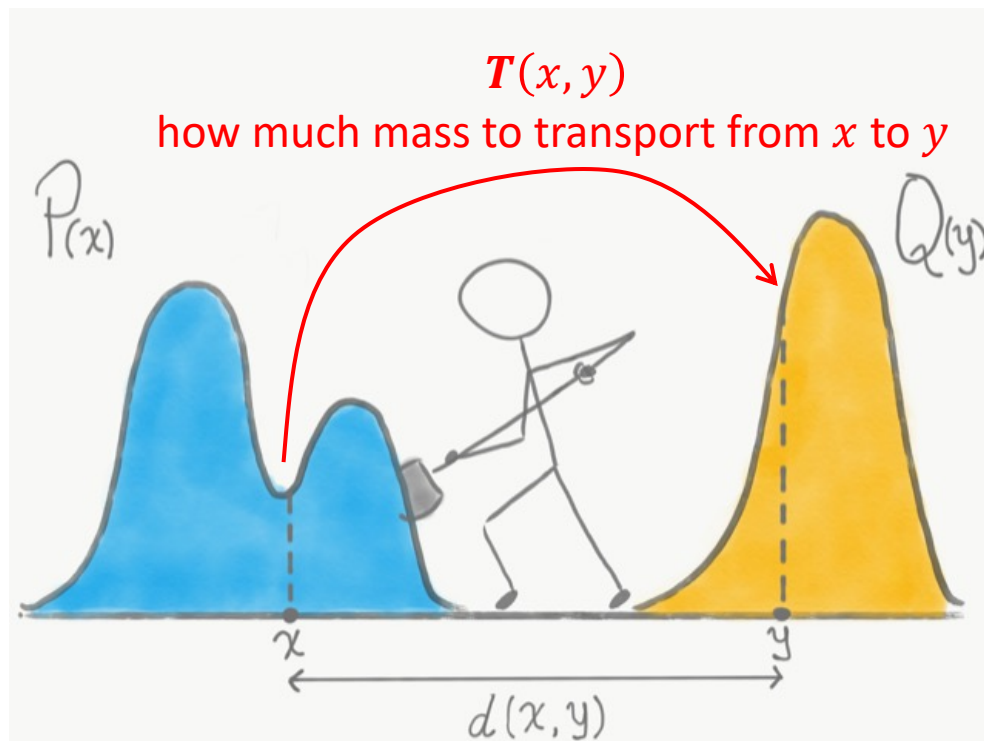
# Preliminaries



- Optimal Transport

- How to transport  $P(x)$  to  $Q(x)$  with minimum effort?

$$\min_T \sum_{x \in P, y \in Q} T(x, y) d(x, y)$$



# Preliminaries

- Wasserstein distance

- Given:

- Two distribution density  $P_1 \in \mathbb{R}^{n_1}, P_2 \in \mathbb{R}^{n_2}$
- Attribute matrices  $X_1 \in \mathbb{R}^{n_1 \times d}, X_2 \in \mathbb{R}^{n_2 \times d}$
- A cross-cost matrix  $C \in \mathbb{R}^{n_1 \times n_2}$  based on  $X_1$  and  $X_2$

- Output:

- The p-Wasserstein distance between  $P_1$  and  $P_2$ :

$$W_p(P_1, P_2) = \left( \min_{T \in \Pi(P_1, P_2)} \sum_{\substack{x \in P_1 \\ y \in P_2}} C^p(x, y) T(x, y) \right)^{1/p}$$

$$\begin{cases} \sum_{x \in P_1} T(x, y) = P_2 \\ \sum_{y \in P_2} T(x, y) = P_1 \end{cases}$$

Minimum effort of transporting  $P_1$  to  $P_2$  in terms of distance between samples



# Preliminaries



- Gromov-Wasserstein distance

- Given:

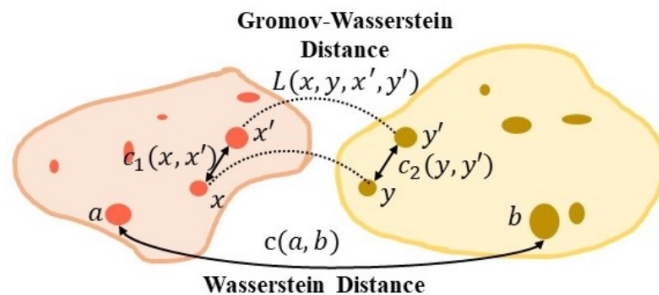
- Two distribution density  $P_1 \in \mathbb{R}^{n_1}, P_2 \in \mathbb{R}^{n_2}$
- Two intra-cost matrices  $C_{P_1} \in \mathbb{R}^{n_1 \times n_1}, C_{P_2} \in \mathbb{R}^{n_2 \times n_2}$

- Output:

- The p-GW distance between  $P_1$  and  $P_2$ :

$$GW_p(P_1, P_2) = \min_{T \in \Pi(P_1, P_2)} \left( \sum_{\substack{x_i, x_j \in P_1 \\ y_l, y_m \in P_2}} \left( C_{P_1}(x_i, x_j) - C_{P_2}(y_l, y_m) \right)^p T(x_i, y_l) T(x_j, y_m) \right)^{1/p}$$

$\downarrow$   
 distance between sample pairs;  
 Similar pairwise relation across distributions



Minimum effort of transporting  $P_1$  to  $P_2$  in terms of distance between sample pairs

# Preliminaries

- Fused Gromov-Wasserstein (FGW) distance
  - Linear combination of  $W_p$  and  $GW_p$

$$\begin{aligned} & \text{FGW}_p(\mathbf{P}_1, \mathbf{P}_2) \\ &= \min_{T \in \Pi(\mathbf{P}_1, \mathbf{P}_2)} \left[ \sum_{\substack{x_i, x_j \in \mathbf{P}_1 \\ y_l, y_m \in \mathbf{P}_2}} \left( (1 - \alpha) \mathbf{C}^p(x_i, y_l) + \alpha \left( \mathbf{C}_{\mathbf{P}_1}(x_i, x_j) - \mathbf{C}_{\mathbf{P}_2}(y_l, y_m) \right)^p \right) \mathbf{T}(x_i, y_l) \mathbf{T}(x_j, y_m) \right]^{1/p} \end{aligned}$$

- Special case:  $L_2$  norm as cross-cost,  $p=2$

- $W_2^2(\mathbf{P}_1, \mathbf{P}_2) = \min_{T \in \Pi(\mathbf{P}_1, \mathbf{P}_2)} \langle \mathbf{C}, \mathbf{T} \rangle; \mathbf{C}(x, y) = \|\mathbf{X}_1(x) - \mathbf{X}_2(y)\|_2^2$
- $GW_2^2(\mathbf{P}_1, \mathbf{P}_2) = \min_{T \in \Pi(\mathbf{P}_1, \mathbf{P}_2)} \langle \mathbf{L}, \mathbf{T} \rangle; \mathbf{L} = \mathbf{C}_{\mathbf{P}_1}^2 \mathbf{P}_1 \mathbf{1}_{n_2}^T + \mathbf{1}_{n_1} \mathbf{P}_2^T \mathbf{C}_{\mathbf{P}_2}^2 - 2 \mathbf{C}_{\mathbf{P}_1} \mathbf{T} \mathbf{C}_{\mathbf{P}_2}^T$
- $\text{FGW}_2^2(\mathbf{P}_1, \mathbf{P}_2) = \min_{T \in \Pi(\mathbf{P}_1, \mathbf{P}_2)} \langle (1 - \alpha) \mathbf{C} + \alpha \mathbf{L}, \mathbf{T} \rangle$

# Outline



- Motivation
- Preliminaries
- Fused Gromov-Wasserstein barycenter
  - Problem definition
  - Optimization
- Experiments
- Takeaways

# FGW barycenter

- A distribution close to given distributions in terms of FGW distance
- Given:

- distribution density  $P_1, \dots, P_K$
- weight for each distribution  $\lambda_1, \dots, \lambda_K$
- intra-cost matrices  $C_{P_1}, \dots, C_{P_K}$
- Barycenter density  $Q \in \mathbb{R}^m$

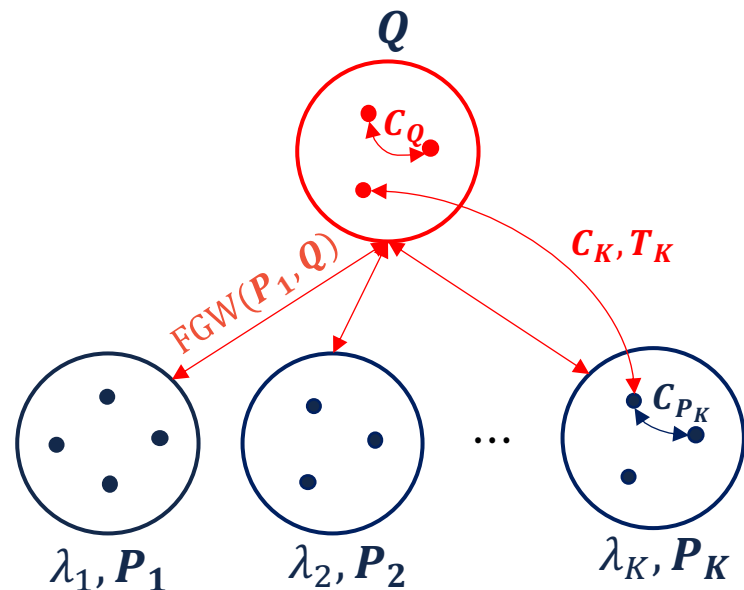
- Output:

- Intra-cost matrix  $C_Q \in \mathbb{R}^{m \times m}$
- Attribute matrix  $X_Q \in \mathbb{R}^{m \times d}$

$$\operatorname{argmin}_{C_Q, X_Q} \sum_{k=1}^K \lambda_k \operatorname{FGW}(P_k, Q)$$

$$= \operatorname{argmin}_{C_Q, X_Q} \sum_{k=1}^K \lambda_k \min_{T_k \in \Pi(P_k, Q)} \langle (1 - \alpha) C_k + \alpha L_k, T_k \rangle$$

$$\text{where} \begin{cases} C_k(x, y) = \|X_{P_k}(x, :) - X_Q(y, :)\|_2 \\ L_k = C_{P_k}^2 P_k \mathbf{1}_m^T + \mathbf{1}_{n_k} Q^T C_Q^{2^T} - 2 C_{P_k} T_k C_Q^T \end{cases}$$



# FGW barycenter



$$\operatorname{argmin}_{\mathbf{C}_Q, \mathbf{X}_Q} \sum_{k=1}^K \lambda_k \min_{\mathbf{T}_k \in \Pi(\mathbf{P}_k, \mathbf{Q})} \langle (1 - \alpha) \mathbf{C}_k + \alpha \mathbf{L}_k, \mathbf{T}_k \rangle$$

where  $\begin{cases} \mathbf{C}_k(x, y) = \|\mathbf{X}_{\mathbf{P}_k}(x, :) - \mathbf{X}_Q(y, :)\|_2 \\ \mathbf{L}_k = \mathbf{C}_{\mathbf{P}_k}^2 \mathbf{P}_k \mathbf{1}_m^T + \mathbf{1}_{n_k} \mathbf{Q}^T \mathbf{C}_Q^{2^T} - 2 \mathbf{C}_{\mathbf{P}_k} \mathbf{T}_k \mathbf{C}_Q^T \end{cases}$

- Block Coordinate Descent
  - Iteratively minimize w.r.t.  $\mathbf{T}_k, \mathbf{X}_Q, \mathbf{C}_Q$
- Fix  $\mathbf{X}_Q^{(t-1)}, \mathbf{C}_Q^{(t-1)}$ ; Optimize  $\mathbf{T}_k^{(t)}$ 
  - *Observation:  $\mathbf{T}_1^{(t)}, \dots, \mathbf{T}_K^{(t)}$  are decoupled*
  - Solve  $K$  problems respectively

$$\mathbf{T}_k^{(t)} = \operatorname{argmin}_{\mathbf{T} \in \Pi(\mathbf{P}_k, \mathbf{Q})} \langle (1 - \alpha) \mathbf{C}_k^{(t-1)} + \alpha \mathbf{L}_k^{(t-1)}, \mathbf{T} \rangle$$

- Sinkhorn algorithm for solution
  - Entropy regularization  $\rightarrow$  strict convexity
  - Iterative matrix scaling  $\rightarrow$  coupling constraint

[4] Vayer, Titouan, et al. "Fused gromov-wasserstein distance for structured objects." Algorithms 13.9 (2020): 212.

[5] Cuturi, Marco. "Sinkhorn distances: Lightspeed computation of optimal transport." Advances in neural information processing systems 26 (2013).

# FGW barycenter



- Fix  $\mathbf{T}_k^{(t)}, \mathbf{X}_Q^{(t-1)}$ ; Optimize  $\mathbf{C}_Q^{(t)}$

$$\operatorname{argmin}_{\mathbf{C}_Q, \mathbf{X}_Q} \sum_{k=1}^K \lambda_k \min_{\mathbf{T}_k \in \Pi(\mathbf{P}_k, \mathbf{Q})} \langle (1 - \alpha) \mathbf{C}_k + \alpha \mathbf{L}_k, \mathbf{T}_k \rangle$$

$$\text{where} \begin{cases} \mathbf{C}_k(x, y) = \|\mathbf{X}_{\mathbf{P}_k}(x, :) - \mathbf{X}_Q(y, :)\|_2 \\ \mathbf{L}_k = \mathbf{C}_{\mathbf{P}_k}^2 \mathbf{P}_k \mathbf{1}_m^T + \mathbf{1}_{n_k} \mathbf{Q}^T \mathbf{C}_Q^{2T} - 2 \mathbf{C}_{\mathbf{P}_k} \mathbf{T}_k \mathbf{C}_Q^T \end{cases}$$

$$\mathbf{C}_Q^{(t)} = \operatorname{argmin}_{\mathbf{C}_Q} \sum_{k=1}^K \lambda_k \left\langle \mathbf{1}_{n_k} \mathbf{Q}^T \mathbf{C}_Q^{2T} - 2 \mathbf{C}_{\mathbf{P}_k} \mathbf{T}_k^{(t)} \mathbf{C}_Q^T, \mathbf{T}_k^{(t)} \right\rangle$$

– First-order optimality

$$\mathbf{C}_Q^{(t)} = \frac{\sum_{k=1}^K \lambda_k \mathbf{T}_k^{(t)T} \mathbf{C}_{\mathbf{P}_k} \mathbf{T}_k^{(t)}}{\mathbf{Q} \mathbf{Q}^T}$$

Average of  $\mathbf{C}_{\mathbf{P}_k}$  based on  $\lambda_k$  and  $\mathbf{T}_k^{(t)}$

$$\left( \mathbf{T}_k^{(t)T} \otimes \mathbf{T}_k^{(t)T} \right) \operatorname{vec}(\mathbf{C}_{\mathbf{P}_k})$$

row  $(i, j)$ , column  $(l, m)$ :  $\mathbf{T}_k^{(t)}(l, i) \mathbf{T}_k^{(t)}(m, j)$   
 Similarity between intra-relation  $(l, m) \in \mathbf{P}_k$   
 and  $(i, j) \in \mathbf{Q}$

# FGW barycenter



- Fix  $T_k^{(t)}, C_Q^{(t)}$ ; Optimize  $X_Q^{(t)}$

$$\operatorname{argmin}_{C_Q, X_Q} \sum_{k=1}^K \lambda_k \min_{T_k \in \Pi(P_K, Q)} \langle (1 - \alpha) C_k + \alpha L_k, T_k \rangle$$

where  $\begin{cases} C_k(x, y) = \|X_{P_k}(x, :) - X_Q(y, :)\|_2 \\ L_k = C_{P_k}^2 P_k \mathbf{1}_m^T + \mathbf{1}_{n_k} Q^T C_Q^{2^T} - 2 C_{P_k} T_k C_Q^T \end{cases}$

$$X_Q^{(t)} = \operatorname{argmin}_{X_Q} \sum_{k=1}^K \lambda_k \langle C_k^{(t)}, T_k^{(t)} \rangle$$

$$\text{where } C_k^{(t)} = \underbrace{\operatorname{diag}(X_k X_k^T) \mathbf{1}_m^T}_{\text{constant}} + \underbrace{\mathbf{1}_{n_k}}_{\mathbf{1}_{n_k} T_k = Q} \operatorname{diag}(X_Q^{(t)} X_Q^{(t)T}) - 2 X_k X_Q^{(t)T}$$

$$X_Q^{(t)} = \operatorname{argmin}_{X_Q} \sum_{k=1}^K \lambda_k \left\| \operatorname{diag}(Q^{\frac{1}{2}}) X_Q - \operatorname{diag}(Q^{-\frac{1}{2}}) T_k^{(t)T} X_k \right\|^2$$

$$X_Q^{(t)} = \operatorname{diag}(Q^{-1}) \sum_{k=1}^K \lambda_k T_k^{(t)T} X_k$$

Average of node attribute  $X_k$  based on  $\lambda_k$  and  $T_k^{(t)}$

# FGW barycenter



---

**Algorithm 1** FGW barycenter

---

**Input** distributions  $P_1, \dots, P_K$ ; weight  $\lambda_1, \dots, \lambda_K$ ; intra-cost matrices  $C_{P_1}, \dots, C_{P_K}$ ; attribute matrices  $X_1, \dots, X_K$ ; barycenter distribution  $Q$

**Output** intra-cost matrix  $C_Q$ ; attribute matrix  $X_Q$

- 1: Initialize  $T_k^{(0)}, C_Q^{(0)}, X_Q^{(0)}$ ;
  - 2: **for**  $t = 1, 2, \dots$  **do**
  - 3:   Update  $T_k^{(t)}$  based on  $C_Q^{(t-1)}, X_Q^{(t-1)}$  via Sinkhorn algorithm.
  - 4:   Update intra-cost  $C_Q^{(t)} = \frac{\sum_{k=1}^K \lambda_k T_k^{(t)T} C_{P_k} T_k^{(t)}}{Q Q^T}$ .
  - 5:   Update attribute matrix  $X_Q^{(t)} = \text{diag}(Q^{-1}) \sum_{k=1}^K \lambda_k T_k^{(t)T} X_k$ .
  - 6: **end for**
  - 7: **return**  $T_k, C_Q, X_Q$ .
- 

[2] Peyré, Gabriel, Marco Cuturi, and Justin Solomon. "Gromov-wasserstein averaging of kernel and distance matrices." International Conference on Machine Learning. PMLR, 2016.

[4] Vayer, Titouan, et al. "Fused gromov-wasserstein distance for structured objects." Algorithms 13.9 (2020): 212.

[5] Cuturi, Marco. "Sinkhorn distances: Lightspeed computation of optimal transport." Advances in neural information processing systems 26 (2013).



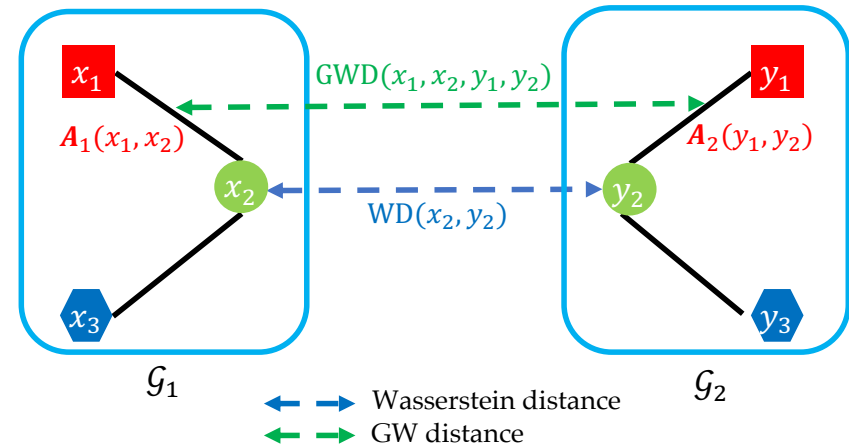
# Outline



- Motivation
- Preliminaries
- Fused Gromov-Wasserstein barycenter
- **Experiments**
- Takeaways

# Experiments

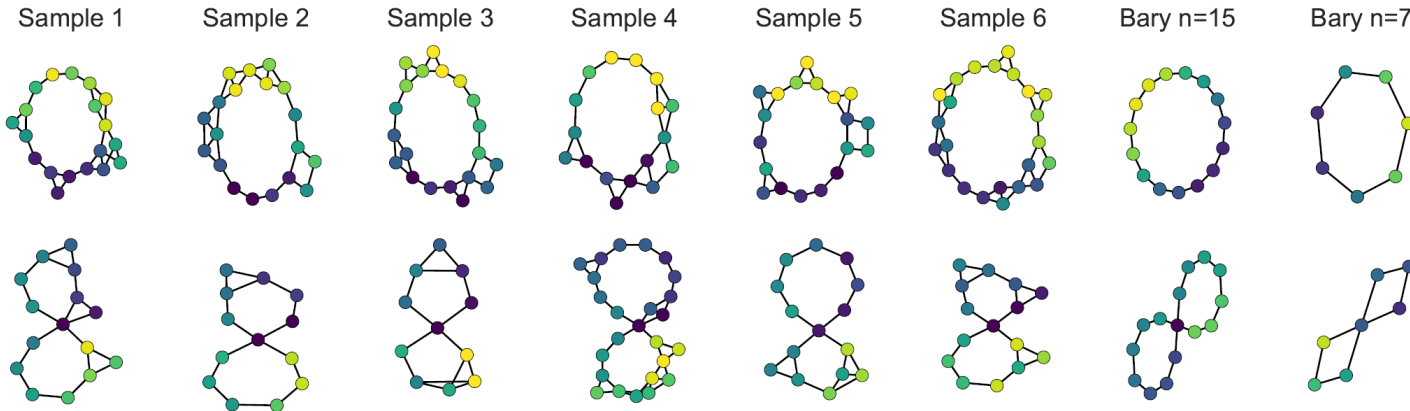
- FGW barycenter for graphs
  - $G_1 = \{V_1, A_1, X_1\}, G_2 = \{V_2, A_2, X_2\}$
  - $P_1 = \frac{1_{|V_1|}}{|V_1|}, P_2 = \frac{1_{|V_2|}}{|V_2|}$
  - Cross-cost  $C$ :  $L_2$  norm between  $X_1, X_2$
  - Intra-cost  $C_{P_1}, C_{P_2}$ : adjacency matrices  $A_1, A_2$
  - WD: node relation; attribute
  - GWD: edge relation; structure



# Experiments



- Graph barycenter



- Graph classification (FGW distance as SVM kernel)

**Table 1.** Average classification accuracy on the graph datasets with vector attributes.

Vector Attributes	BZR	COX2	CUNEIFORM	ENZYMES	PROTEIN	SYNTHETIC
FGW sp	<b>85.12 ± 4.15 *</b>	77.23 ± 4.86	<b>76.67 ± 7.04</b>	71.00 ± 6.76	<b>74.55 ± 2.74</b>	<b>100.00 ± 0.00</b>
HOPPERK	84.15 ± 5.26	<b>79.57 ± 3.46</b>	32.59 ± 8.73	45.33 ± 4.00	71.96 ± 3.22	90.67 ± 4.67
PROPAK	79.51 ± 5.02	77.66 ± 3.95	12.59 ± 6.67	<b>71.67 ± 5.63 *</b>	61.34 ± 4.38	64.67 ± 6.70
PSCN k = 10	80.00 ± 4.47	71.70 ± 3.57	25.19 ± 7.73	26.67 ± 4.77	67.95 ± 11.28	<b>100.00 ± 0.00</b>
PSCN k = 5	82.20 ± 4.23	71.91 ± 3.40	24.81 ± 7.23	27.33 ± 4.16	71.79 ± 3.39	<b>100.00 ± 0.00</b>

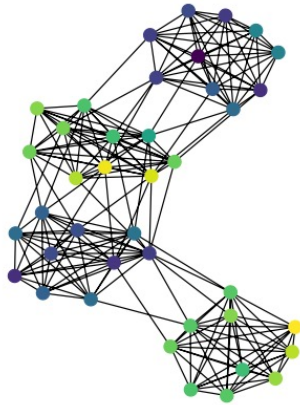
[4] Vayer, Titouan, et al. "Fused gromov-wasserstein distance for structured objects." Algorithms 13.9 (2020): 212.

# Experiments

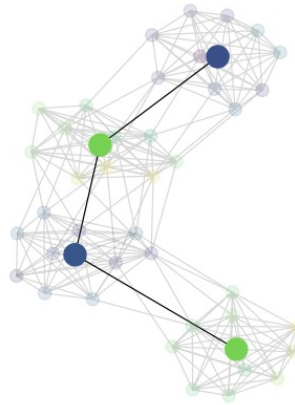


- Clustering

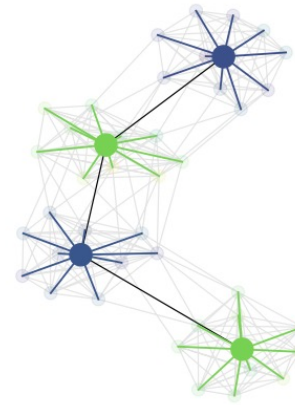
Graph with communities



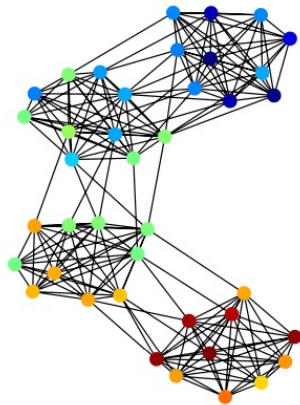
Approximate Graph



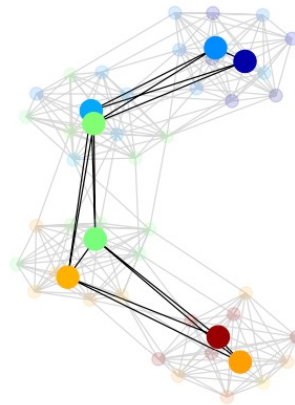
Clustering with transport matrix



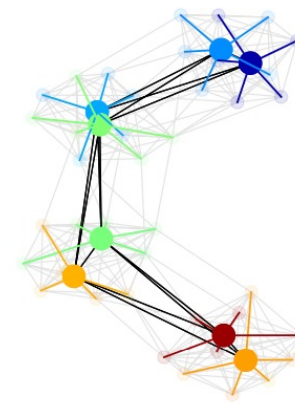
Graph with bimodal communities



Approximate Graph



Clustering with transport matrix



[4] Vayer, Titouan, et al. "Fused gromov-wasserstein distance for structured objects." Algorithms 13.9 (2020): 212.

# Outline



- Motivation
- Preliminaries
- Fused Gromov-Wasserstein barycenter
- Experiments
- Takeaways

# Takeaways

- Motivation
- Preliminaries
  - OT, WD, GWD, FGWD
- Fused Gromov-Wasserstein barycenter
  - Block coordinate descent
  - Optimize transport plan  $\mathbf{T}_k$ : solve K optimal transport problems
  - Optimize intra-cost  $\mathbf{C}_Q$ : Average of  $\mathbf{C}_k$  based on  $\lambda_k$  and  $\mathbf{T}_k$
  - Optimize attribute  $\mathbf{X}_Q$ : Average of  $\mathbf{X}_k$  based on  $\lambda_k$  and  $\mathbf{T}_k$
- Experiments
  - Graph barycenter
  - Graph classification
  - Graph clustering

# References



- [1] Peyré, Gabriel, and Marco Cuturi. "Computational Optimal Transport." arXiv preprint arXiv:1803.00567 (2018).
- [2] Peyré, Gabriel, Marco Cuturi, and Justin Solomon. "Gromov-wasserstein averaging of kernel and distance matrices." International Conference on Machine Learning. PMLR, 2016.
- [3] Chen, Liqun, et al. "Graph optimal transport for cross-domain alignment." International Conference on Machine Learning. PMLR, 2020.
- [4] Vayer, Titouan, et al. "Fused gromov-wasserstein distance for structured objects." Algorithms 13.9 (2020): 212.
- [5] Cuturi, Marco. "Sinkhorn distances: Lightspeed computation of optimal transport." Advances in neural information processing systems 26 (2013).

# Thanks for listening!

## Q&A

Zhichen Zeng

10/13

A black and white photograph of the dome of the University of Illinois State Capitol building. The dome is a large, circular structure with a series of ribs radiating from the top. The top of the dome is decorated with a series of small, ornate finials. The image is taken from a low angle, looking up at the dome.

**I** ILLINOIS