

# Online Graph Dictionary Learning

Ruike Zhu

2/9 2023



UNIVERSITY OF  
**ILLINOIS**  
URBANA-CHAMPAIGN

# Outline

- Introduction
- Preliminaries
- Online graph dictionary learning
- Experiments
- Conclusion



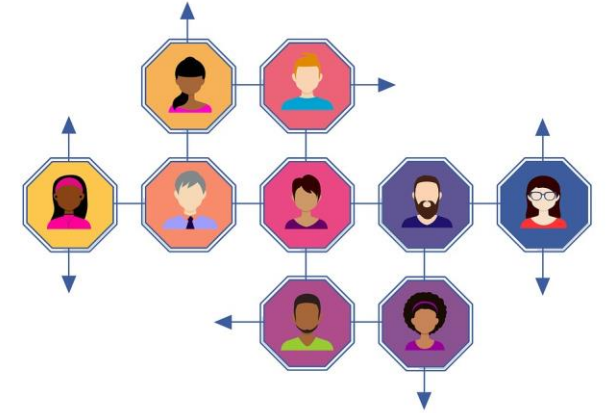
# Outline

- Introduction
- Preliminaries
- Online graph dictionary learning
- Experiments
- Conclusion



# Introduction

- Graphs has been of great interest in the last decades
  - ◇ molecule compounds
  - ◇ brain connectivity
  - ◇ Social networks
  - ◇ ...
- Designing good representations for these data is challenging
  - ◇ non-vectorial
  - ◇ requires dedicated modelling of their representing structures
- Introduce an unsupervised representation learning algorithm



# Outline

- Introduction
- Preliminaries
  - ◇ Dictionary Learning
  - ◇ (Fused) Gromov-Wasserstein for structured data
- Online graph dictionary learning
- Experiments
- Conclusion



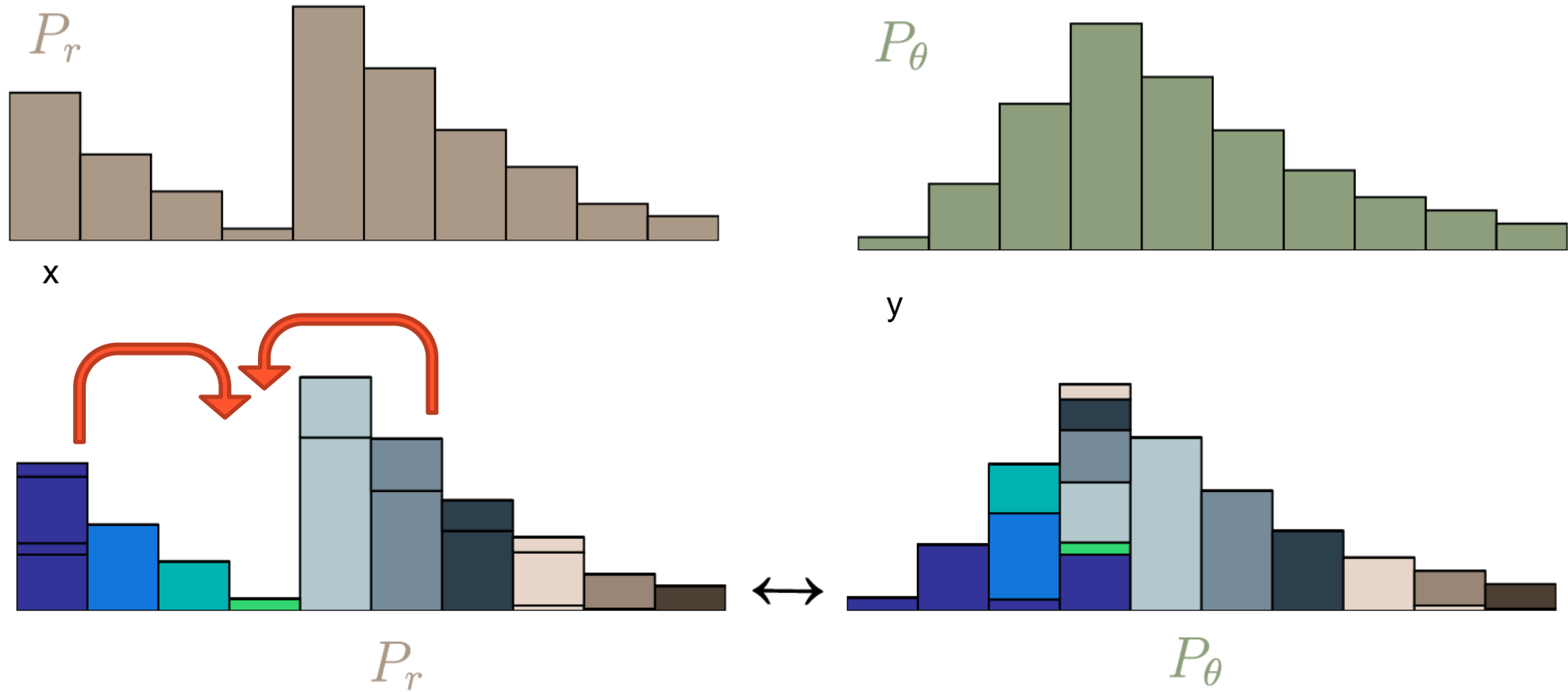
# Preliminaries

- Dictionary learning
  - ◇ A field of unsupervised learning that aims at estimating a linear representation of the data. The linear subspace is defined by the span of a family of vectors, called **atoms**, which constitute a **dictionary**
  - ◇ These atoms are inferred from the input data by minimizing a reconstruction error
- Linear modeling of graph:
  - ◇ Given: a graph  $G = (\mathbf{C}, \mathbf{h})$  with intra-cost matrices  $\mathbf{C} \in \mathbb{R}^{N \times N}$ , node weight  $\mathbf{h} \in \mathbb{R}^N$  and a dictionary  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$
  - ◇ Output: a linear representation  $\sum_{s \in [S]} \omega_s \bar{\mathbf{c}}_s$  of graph  $G$  with minimized reconstruction error

# Preliminaries

$$\begin{aligned}\sum_x T(x, y) &= P_r(y), \\ \sum_y T(x, y) &= P_\theta(x)\end{aligned}$$

- Earth Mover's Distance:  $EMD(P_r, P_\theta) = \min \sum_{x,y} \|x - y\| T(x, y)$



# Preliminaries

- (Fused) Gromov-Wasserstein for graph data

- ◇ Given two graphs  $G^X = (\mathbf{C}^X, \mathbf{h}^X)$ ,  $G^Y = (\mathbf{C}^Y, \mathbf{h}^Y)$

- $\mathbf{C}^X \in \mathbb{R}^{N^X \times N^X}$  and  $\mathbf{C}^Y \in \mathbb{R}^{N^Y \times N^Y}$  are the intra-cost matrices.

- $\mathbf{h}^X, \mathbf{h}^Y$  is a vector of node weights in graph, usually define  $\mathbf{h}^X = \frac{1}{N^X} \mathbf{1}_{N^X}$  if without any prior knowledge

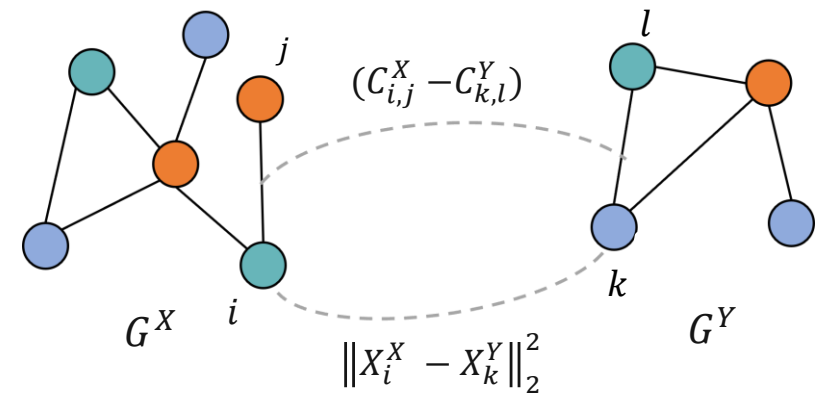
Distribution density

- ◇ The  $GW_p$  distance between  $G^X$  and  $G^Y$  is defined as:

$$\min_{T \in \Pi(\mathbf{h}^X, \mathbf{h}^Y)} \left( \sum_{\substack{i,j \in G^X \\ k,l \in G^Y}} (\mathbf{C}_{i,j}^X - \mathbf{C}_{k,l}^Y)^p T_{ik} T_{jl} \right)$$

- ◇  $\Pi(\mathbf{h}^X, \mathbf{h}^Y) := \{T \in \mathbb{R}^{N^X \times N^Y} | T\mathbf{1}_{N^Y} = \mathbf{h}^X, T\mathbf{1}_{N^X} = \mathbf{h}^Y\}$

- ◇ FGW distance:  $\min_{T \in \Pi(\mathbf{h}^X, \mathbf{h}^Y)} \left( \sum_{\substack{i,j \in G^X \\ k,l \in G^Y}} ((1 - \alpha) \|X_i^X - X_k^Y\|_2^2 + \alpha (\mathbf{C}_{i,j}^X - \mathbf{C}_{k,l}^Y)^p) T_{ik} T_{jl} \right)$



T act as a probabilistic matching of nodes which tends to associate pairs of nodes that have similar pairwise relations in  $\mathbf{C}^X$  and  $\mathbf{C}^Y$



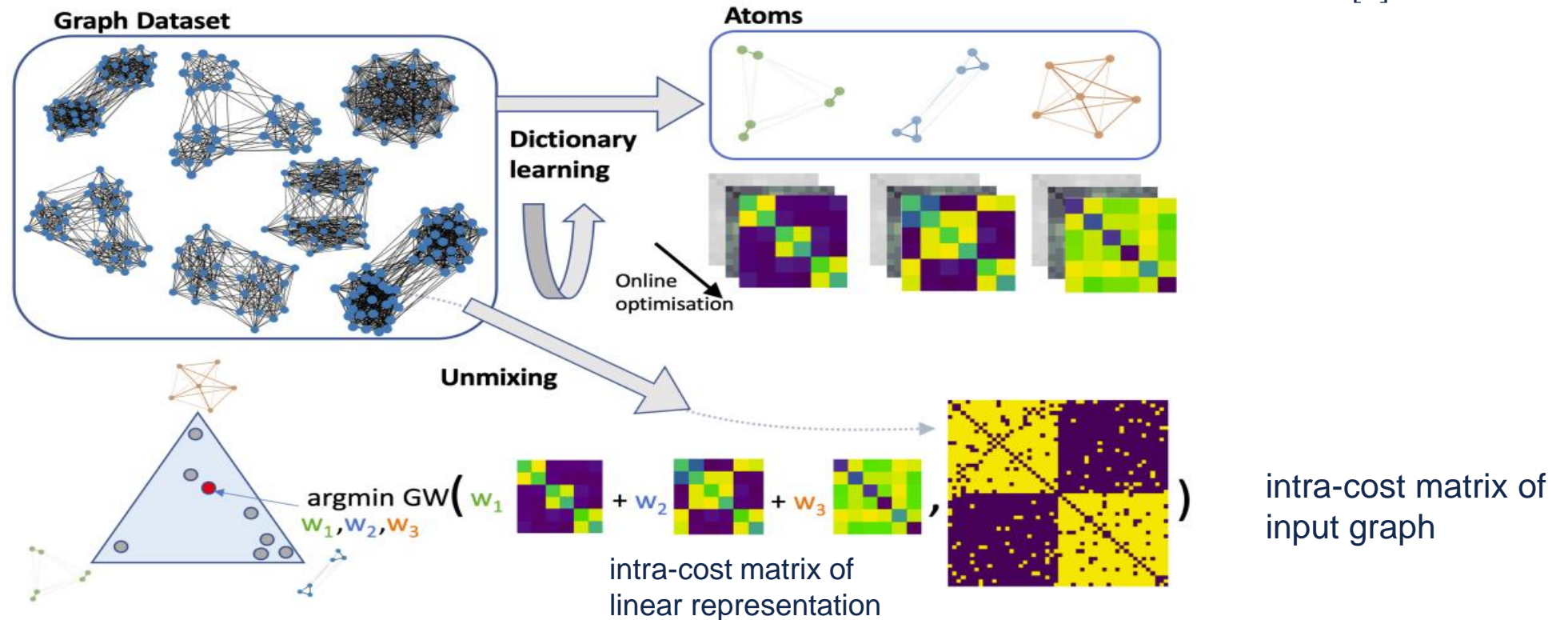
# Outline

- Introduction
- Preliminaries
- Online graph dictionary learning
  - ◇ GW unmixing
  - ◇ Dictionary learning and online algorithm
- Experiments
- Conclusion



# Algorithm overview

- Minimization of the GW distance between the intra-cost matrix associated to the input graph and its linear representation in the dictionary
- - Given a set of graph  $\{G^{(k)}: (\mathbf{C}^{(k)}, \mathbf{h}^{(k)})\}_{k \in [K]}$ 
  - Output the optimal embedding vector  $\omega \in \Sigma_S$  of each graph  $G$  and a dictionary  $\{\bar{\mathbf{C}}_s\}_{s \in [S]}$



# Algorithm overview

- - Given a set of graph  $\{G^{(k)}: (\mathbf{c}^{(k)}, \mathbf{h}^{(k)})\}_{k \in [K]}$ 
  - Output the optimal embedding vector  $\boldsymbol{\omega} \in \Sigma_S$  of each graph  $G$  and a dictionary  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$
- Workflow of algorithm:

```
1 Randomly initialize  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$ 
2 For each batch of graph:
3     Initialize  $\boldsymbol{\omega} = \frac{1}{S} \mathbf{1}_S$ 
4     repeat
5         Fix dictionary  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$  and embedding vector  $\boldsymbol{\omega}$ , update transport plan  $\mathbf{T}$ 
6         Fix dictionary  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$  and transport plan  $\mathbf{T}$ , update embedding vector  $\boldsymbol{\omega}$ 
7     until converge
8 Fix  $\boldsymbol{\omega}$  and  $\mathbf{T}$ , update  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$ , projected gradient step
```

GW  
Unmixing

# GW Unmixing

- Minimization of the GW distance between the similarity matrix associated to the graph and its linear representation in the dictionary
- - Given a graph  $G = (\mathbf{C}, \mathbf{h})$   
a dictionary  $\{\bar{\mathbf{C}}_s\}_{s \in [S]}$  with  $S$  atoms of which each atom contains  $N_s$  nodes
- Output the optimal transport matrix  $\mathbf{T}$ ,  
embedding vector  $\boldsymbol{\omega}$  of graph  $G$
- Object function of unmixing problem given dictionary:

$$\diamond \quad \min_{\boldsymbol{\omega} \in \Sigma_S} \left( GW_2^2 \left( \mathbf{C}, \sum_{s \in [S]} \omega_s \bar{\mathbf{C}}_s \right) - \boxed{\lambda} \|\boldsymbol{\omega}\|_2^2 \right)$$

$\lambda$  is the negative quadratic regularization promoting sparsity [1]

When  $\mathbf{C}_s$  is adjacency matrix, this represent the probabilities of connection between the nodes in linear representation

**Fix  $\{\bar{\mathbf{C}}_s\}_{s \in [S]}$ ,  $\omega$  – update  $T$**

$$\diamond \min_{\omega \in \Sigma_S} \left( GW_2^2 \left( \mathbf{C}, \sum_{s \in [S]} \omega_s \bar{\mathbf{C}}_s \right) - \lambda \|\omega\|_2^2 \right) \quad (1)$$

---

**Algorithm 3** BCD for GW unmixing problem

---

- 1: Initialize  $\mathbf{w} = \frac{1}{S} \mathbf{1}_S$
  - 2: **repeat**
  - 3:   Compute OT matrix  $T$  of  $GW_2^2 \left( \mathbf{C}, \tilde{\mathbf{C}}(\mathbf{w}) \right)$ , with CG algorithm (Vayer et al., 2018, Alg.1 & 2).
  - 4:   Compute the optimal  $\mathbf{w}$  solving equation 1 for a fixed  $T$  with CG algorithm 4
  - 5: **until** convergence
- 

- $\diamond GW_2^2 \left( \mathbf{C}, \sum_{s \in [S]} \omega_s \bar{\mathbf{C}}_s \right) = \min_{T \in \Pi(\mathbf{h}, \bar{\mathbf{h}}_S)} \langle L, T \rangle; L = \mathbf{C}^2 \mathbf{h} \mathbf{1}_N^T + \mathbf{1}_{N_S} \mathbf{h}_S^T (\tilde{\mathbf{C}}(\omega)^2)^T - 2 \mathbf{C} T \tilde{\mathbf{C}}(\omega)^T$  according to equation (5) in [3]
- $\diamond T^{(t)} = \operatorname{argmin}_{T \in \Pi(\mathbf{h}, \bar{\mathbf{h}}_S)} \langle \mathbf{C}^2 \mathbf{h} \mathbf{1}_N^T + \mathbf{1}_{N_S} \mathbf{h}_S^T (\tilde{\mathbf{C}}(\omega)^2)^T - 2 \mathbf{C} T \tilde{\mathbf{C}}(\omega)^T, T \rangle; \tilde{\mathbf{C}}(\omega) = \sum_{s \in [S]} \omega_s \bar{\mathbf{C}}_s,$
- $\diamond$  This corresponds to an optimal transport problem, can be solved iteratively by the entropic regularization-based method



[1] Vincent-Cuaz, Cédric, et al. "Online graph dictionary learning." International Conference on Machine Learning. PMLR, 2021.

[3] Xu, Hongteng. "Gromov-Wasserstein factorization models for graph clustering." Proceedings of the AAAI conference on artificial intelligence. Vol. 34. No. 04. 2020.

**Fix  $\{\bar{\mathbf{C}}_s\}_{s \in [S]}$ ,  $\mathbf{T}$  – update  $\mathbf{w}$**

$$\varepsilon(\mathbf{C}, \tilde{\mathbf{C}}(\mathbf{w}), \mathbf{T}) = GW_2^2\left(\mathbf{C}, \sum_{s \in [S]} \omega_s \bar{\mathbf{C}}_s\right) - \lambda \|\mathbf{w}\|_2^2$$

---

**Algorithm 4** CG for solving GW unmixing problem *w.r.t*  $\mathbf{w}$  given  $\mathbf{T}$

---

- 1: **repeat**
- 2:   Compute  $\mathbf{g}$ , gradients *w.r.t*  $\mathbf{w}$  of  $\mathcal{E}(\mathbf{C}, \tilde{\mathbf{C}}(\mathbf{w}), \mathbf{T})$  following equation 2
- 3:   Find direction  $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \Sigma_S} \mathbf{x}^T \mathbf{g}$
- 4:   Line-search: denoting  $\mathbf{z}(\gamma) = \gamma \mathbf{x}^* + (1 - \gamma) \mathbf{w}$ , **#ensure the constrain of  $\|\mathbf{w}\|_1 = 1$**

$$\gamma^* = \arg \min_{\gamma \in (0,1)} \mathcal{E}(\mathbf{C}, \tilde{\mathbf{C}}(\mathbf{z}(\gamma)), \mathbf{T}) = \arg \min_{\gamma \in (0,1)} a\gamma^2 + b\gamma + c$$

- 5:    $\mathbf{w} \leftarrow \mathbf{z}(\gamma^*)$
  - 6: **until** convergence
- 

**#  $\varepsilon(\mathbf{C}, \tilde{\mathbf{C}}(\mathbf{w}), \mathbf{T})$  can be converted to second-order polynomial equation with coefficients as below**

Partial derivates of the GW objective  $\varepsilon$

$$\diamond \quad \frac{\partial \varepsilon}{\partial \omega_s}(\mathbf{C}, \tilde{\mathbf{C}}(\mathbf{w}), \mathbf{T}) = 2tr\left\{\left(\bar{\mathbf{C}}_s \odot \tilde{\mathbf{C}}(\mathbf{w})\right) \mathbf{h} \mathbf{h}^T - \bar{\mathbf{C}}_s \mathbf{T}^T \mathbf{C}^T \mathbf{T}\right\} \quad (2)$$

$$a = tr\left\{\left(\tilde{\mathbf{C}}(\mathbf{x}^* - \mathbf{w}) \odot \tilde{\mathbf{C}}(\mathbf{x}^* - \mathbf{w})\right) \mathbf{h} \mathbf{h}^T\right\} - \lambda \|\mathbf{x}^* - \mathbf{w}\|_2^2$$

$$b = 2tr\left\{\left(\tilde{\mathbf{C}}(\mathbf{x}^* - \mathbf{w}) \odot \tilde{\mathbf{C}}(\mathbf{w})\right) \mathbf{h} \mathbf{h}^T - \tilde{\mathbf{C}}(\mathbf{x}^* - \mathbf{w}) \mathbf{T}^T \mathbf{C}^T \mathbf{T}\right\} - 2\lambda \langle \mathbf{w}, \mathbf{x} - \mathbf{w} \rangle$$

# Online algorithm

- Assume the dictionary is unknown and need to be estimated from dataset
- - Given  $K$  graphs  $\{G^{(k)}: (\mathbf{C}^{(k)}, \mathbf{h}^{(k)})\}_{k \in [K]}$ 
  - Output the optimal embedding vector  $\boldsymbol{\omega}$  of each graph  $G$  and a dictionary  $\{\bar{\mathbf{C}}_s\}_{s \in [S]}$
- Object function:

$$\diamond \min_{\{\boldsymbol{\omega}^{(k)}\}_{k \in [K]}, \{\bar{\mathbf{C}}_s\}_{s \in [S]}} \sum_{k=1}^K (GW_2^2(\mathbf{C}^{(k)}, \sum_{s \in [S]} \omega_s^{(k)} \bar{\mathbf{C}}_s) - \lambda \|\boldsymbol{\omega}^{(k)}\|_2^2)$$

$$\diamond \boldsymbol{\omega}^{(k)} \in \Sigma_S, \bar{\mathbf{C}}_s \in S_N(\mathbb{R})$$

# Fix $T, \omega$ – update $\{\bar{\mathbf{C}}_s\}_{s \in [S]}$

$$\min_{\substack{\{\omega^{(k)}\}_{k \in [B]} \\ \{\bar{\mathbf{C}}_s\}_{s \in [S]}}} \sum_{k=1}^B (GW_2^2(\mathbf{C}^{(k)}, \sum_{s \in [S]} \omega_s^{(k)} \bar{\mathbf{C}}_s) - \lambda \|\omega^{(k)}\|_2^2)$$

---

**Algorithm 5** GDL: stochastic update of atoms  $\{\bar{\mathbf{C}}_s\}_{s \in [S]}$

---

- 1: Sample a minibatch of graphs  $\mathcal{B} := \{\mathbf{C}^{(k)}\}_{k \in \mathcal{B}}$ .
- 2: Compute optimal  $\{(\mathbf{w}^{(k)}, \mathbf{T}^{(k)})\}_{k \in [B]}$  by solving B independent unmixing problems with Alg.3.
- 3: Projected gradient step with estimated gradients  $\tilde{\nabla}_{\bar{\mathbf{C}}_s}$  (see equation 54),  $\forall s \in [S]$ :

$$\bar{\mathbf{C}}_s \leftarrow Proj_{S_N(\mathbb{R})}(\bar{\mathbf{C}}_s - \eta_C \tilde{\nabla}_{\bar{\mathbf{C}}_s}) \quad \# \text{ ensure constrain of } \bar{\mathbf{C}}_s, \text{ need to be symmetric matrix}$$


---

Estimated gradients *w.r.t*  $\{\bar{\mathbf{C}}_s\}$  over a minibatch of graphs  $\mathcal{B} := \{\mathbf{C}^{(k)}\}_{k \in \mathcal{B}}$  given unmixing solutions  $\{(\mathbf{w}^{(k)}, \mathbf{T}^{(k)})\}_{k \in [B]}$  read:

$$\tilde{\nabla}_{\bar{\mathbf{C}}_s} \left( \sum_{k \in \mathcal{B}} \mathcal{E}(\mathbf{C}^{(k)}, \tilde{\mathbf{C}}(\mathbf{w}^{(k)}), \mathbf{T}^{(k)}) \right) = \frac{2}{B} \sum_{k \in \mathcal{B}} w_s^{(k)} \{ \tilde{\mathbf{C}}(\mathbf{w}^{(k)}) \odot \mathbf{h} \mathbf{h}^\top - \mathbf{T}^{(k)\top} \mathbf{C}^{(k)\top} \mathbf{T}^{(k)} \} \quad (54)$$



# Algorithm overview

- - Given a set of graph  $\{G^{(k)}: (\mathbf{c}^{(k)}, \mathbf{h}^{(k)})\}_{k \in [K]}$ 
  - Output the optimal embedding vector  $\boldsymbol{\omega} \in \Sigma_S$  of each graph  $G$  and a dictionary  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$
- Workflow of algorithm:

- 1 Randomly initialize  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$
- 2 For each batch of graph:
- 3     Initialize  $\boldsymbol{\omega} = \frac{1}{S} \mathbf{1}_S$
- 4     repeat
- 5         Fix dictionary  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$  and embedding vector  $\boldsymbol{\omega}$ , update transport plan  $\mathbf{T}$
- 6         Fix dictionary  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$  and transport plan  $\mathbf{T}$ , update embedding vector  $\boldsymbol{\omega}$
- 7     until converge
- 8 Fix  $\boldsymbol{\omega}$  and  $\mathbf{T}$ , update  $\{\bar{\mathbf{c}}_s\}_{s \in [S]}$ , projected gradient step

Algorithm 3

Algorithm 4

Algorithm 5

# Outline

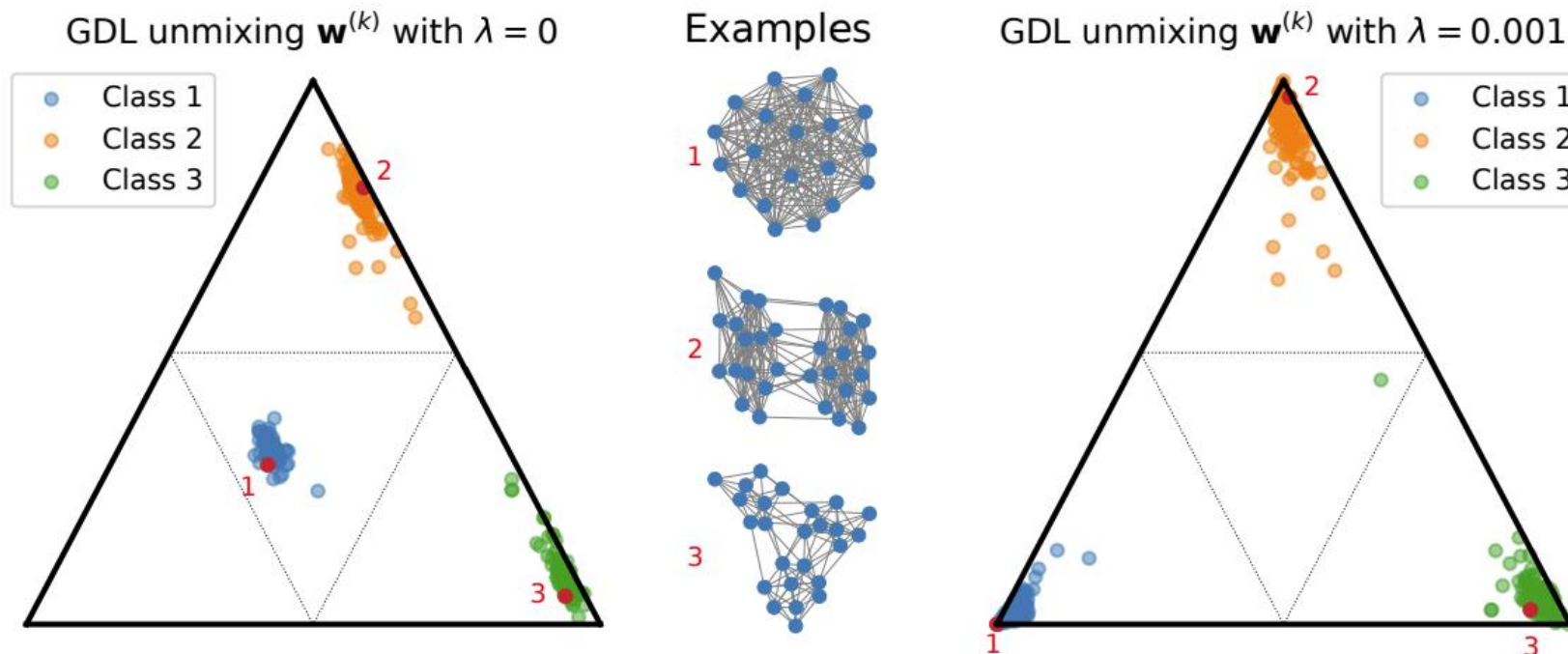
- Introduction
- Preliminaries
- Online graph dictionary learning
- **Experiments**
  - ◇ With synthesis graph
  - ◇ With real-life graph
  - ◇ With Stream graph
- Conclusion



## With synthesis graph

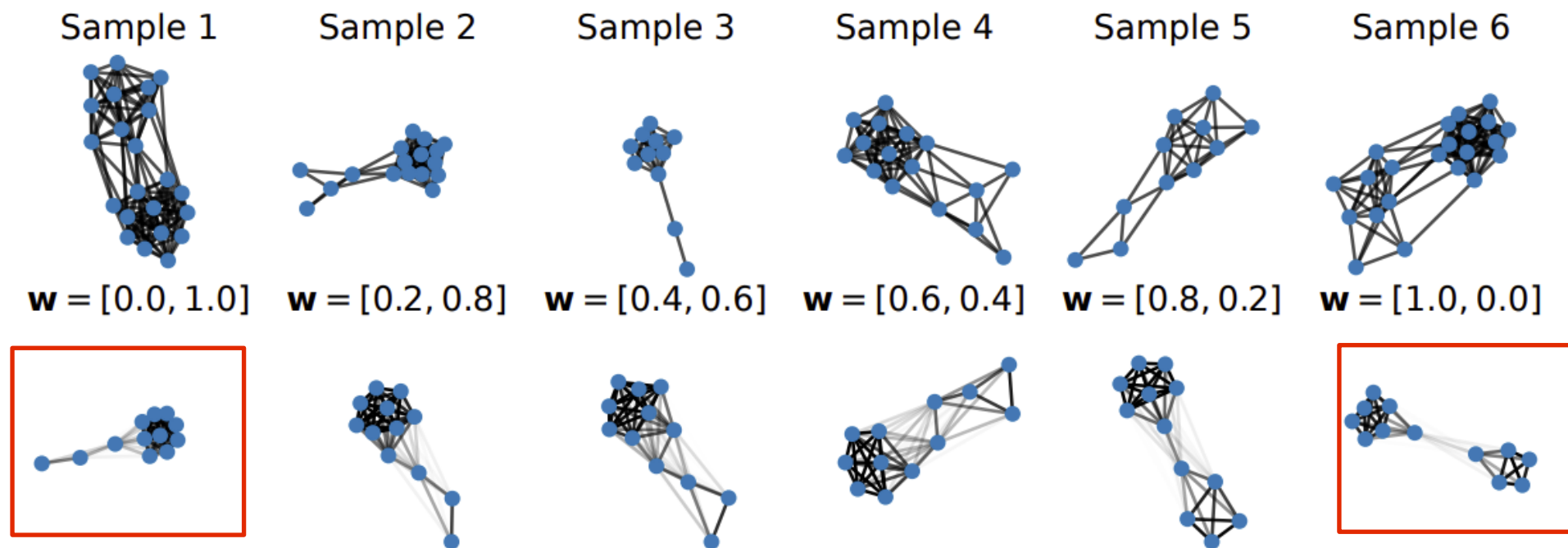
$$\min_{\{\omega^{(k)}\}_{k \in [K]}, \{\bar{\mathbf{c}}_s\}_{s \in [S]}} \sum_{k=1}^K (GW_2^2(\mathbf{c}^{(k)}, \sum_{s \in [S]} \omega_s^{(k)} \bar{\mathbf{c}}_s) - \lambda \|\omega^{(k)}\|_2^2)$$

- Three different generative structures: dense, two clusters and three clusters.
- Nodes are assigned to clusters into equal proportions.
- For each generative structure 100 graphs are sampled.



## With synthesis graph

- Estimate on D2 a dictionary with 2 atoms of order 12. The interpolation between the two estimated atoms for some samples is reported.



On the top, a random sample of real graphs from D2 (two blocks).

On the bottom, reconstructed graphs as linear combination of two estimated atoms (varying proportions for each atom).

# With real-life graph

- Graph clustering

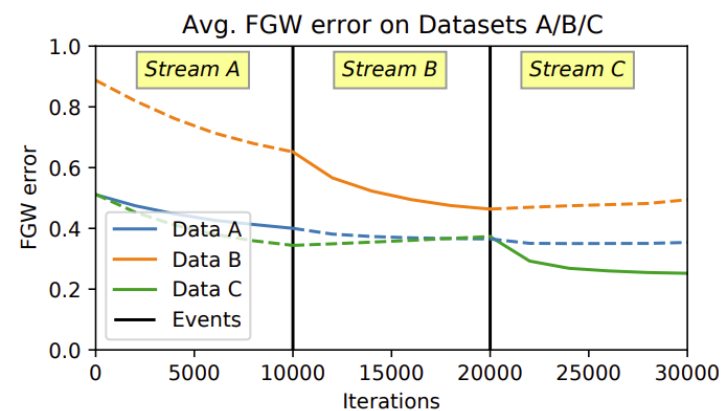
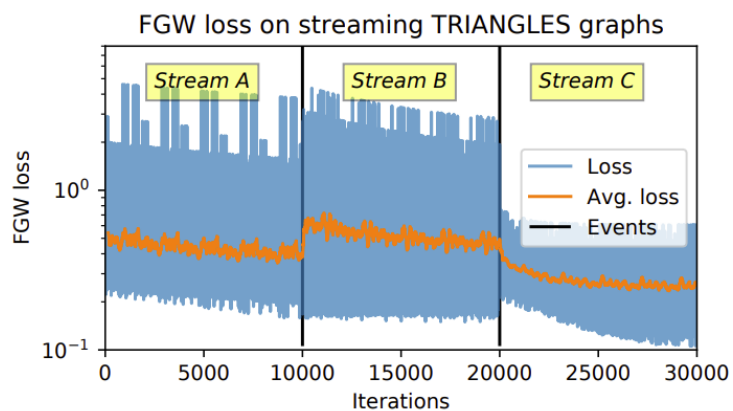
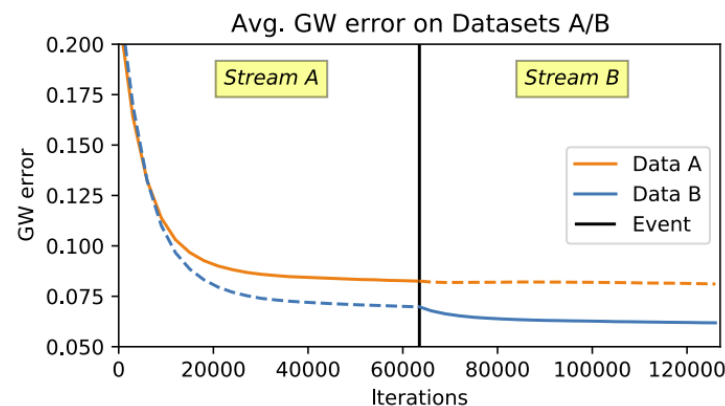
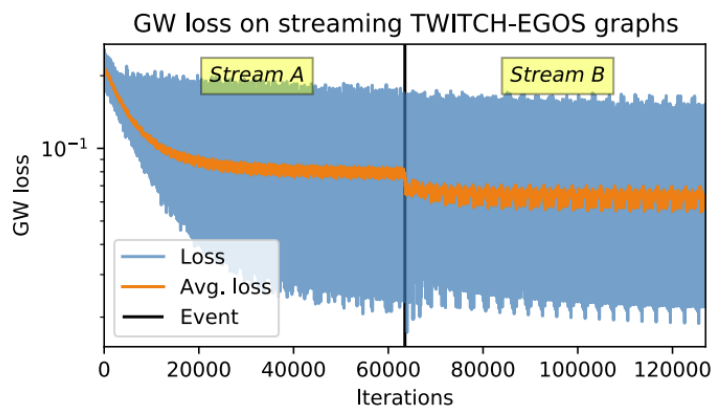
Table 1. Clustering: Rand Index computed for benchmarked approaches on real datasets.

MODELS	NO ATTRIBUTE		DISCRETE ATTRIBUTES		REAL ATTRIBUTES			
	IMDB-B	IMDB-M	MUTAG	PTC-MR	BZR	COX2	ENZYMES	PROTEIN
GDL (ours)	51.32(0.30)	55.08(0.28)	70.02(0.29)	51.53(0.36)	62.59(1.68)	58.39(0.52)	66.97(0.93)	60.22(0.30)
<b>GDL<sub><math>\lambda</math></sub> (ours)</b>	<b>51.64(0.59)</b>	55.41(0.20)	<b>70.89(0.11)</b>	<b>51.90(0.54)</b>	<b>66.42(1.96)</b>	<b>59.48(0.68)</b>	66.79(1.12)	<b>60.49(0.71)</b>
GWF-r	51.24 (0.02)	<b>55.54(0.03)</b>	68.83(1.47)	51.44(0.52)	52.42(2.48)	56.84(0.41)	<b>72.13(0.19)</b>	59.96(0.09)
GWF-f	50.47(0.34)	54.01(0.37)	58.96(1.91)	50.87(0.79)	51.65(2.96)	52.86(0.53)	71.64(0.31)	58.89(0.39)
GW-k	50.32(0.02)	53.65(0.07)	57.56(1.50)	50.44(0.35)	56.72(0.50)	52.48(0.12)	66.33(1.42)	50.08(0.01)
SC	50.11(0.10)	54.40(9.45)	50.82(2.71)	50.45(0.31)	42.73(7.06)	41.32(6.07)	70.74(10.60)	49.92(1.23)

With negative quadratic regularization

# Stream graph

- Dataset: TWITCH-EGOS with 2 atoms of 14 nodes each (top)  
TRIANGLES with 4 atoms of 17 nodes each (bottom).



# Outline

- Introduction
- Preliminaries
- Online graph dictionary learning
- Experiments
- Conclusion



# Conclusion

- Introduction
- Preliminaries
  - ◇ Dictionary Learning
  - ◇ (Fused) Gromov-Wasserstein for structured data
- Online graph dictionary learning
  - ◇ GW unmixing
  - ◇ Dictionary learning and online algorithm
- Experiments
  - ◇ With synthesis graph
  - ◇ With real-life graph
  - ◇ With Stream graph





# Reference

- [1] Vincent-Cuaz, Cédric, et al. "Online graph dictionary learning." International Conference on Machine Learning. PMLR, 2021.
- [2] Peyré, Gabriel, Marco Cuturi, and Justin Solomon. "Gromov-wasserstein averaging of kernel and distance matrices." International conference on machine learning. PMLR, 2016.
- [3] Xu, Hongtengl. "Gromov-Wasserstein factorization models for graph clustering." Proceedings of the AAAI conference on artificial intelligence. Vol. 34. No. 04. 2020.

Thanks!

