

The background of the slide is a light pink color with a repeating pattern of stylized flowers and leaves. The flowers have five petals and a central stamen, while the leaves are elongated and pointed. The pattern is dense and covers the entire background.

Research Progress Report

Reporter: Minjie Cheng

July 20, 2023

Overview

Research Topic: **Neural Network Modeling based on Optimal Transport.**

- ▶ Revisiting Global Pooling via Regularized OT \Rightarrow OT for Pooling
- ▶ A Quasi-Wasserstein Loss for Learning Graph Neural Networks \Rightarrow OT for GNNs

Revisiting Global Pooling via Regularized OT

- ▶ A Global Pooling outputs **the expectation of samples conditioned on different feature dimensions.**

$$f(\mathbf{X}) = (\mathbf{X} \odot \underbrace{\text{diag}^{-1}(\overbrace{\mathbf{P}\mathbf{1}_N}^{\mathbf{p}=[p_d]})}_{\tilde{\mathbf{P}}=[p_{n|d}]})\mathbf{1}_N = \left\|_{d=1}^D \mathbb{E}_{n \sim p_{n|d}}[x_{dn}], \quad (1)$$

- ▶ The pooling is determined by **the "sample-feature" distribution**

- ▶ Mean-pooling: $f(\mathbf{X}) = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad \Rightarrow \quad \mathbf{P} = [\frac{1}{DN}]$
- ▶ Max-pooling: $f(\mathbf{X}) = \left\|_{d=1}^D \max_n \{x_{dn}\}_{n=1}^N \quad \Rightarrow \quad \mathbf{P} \in \{0, \frac{1}{D}\}^{D \times N}$
- ▶ Attention-pooling: $f(\mathbf{X}) = \mathbf{X}\mathbf{a}_X \quad \Rightarrow \quad \mathbf{P} = \frac{1}{D} \mathbf{1}_D \mathbf{a}_X^T$

Revisiting Global Pooling via Regularized OT

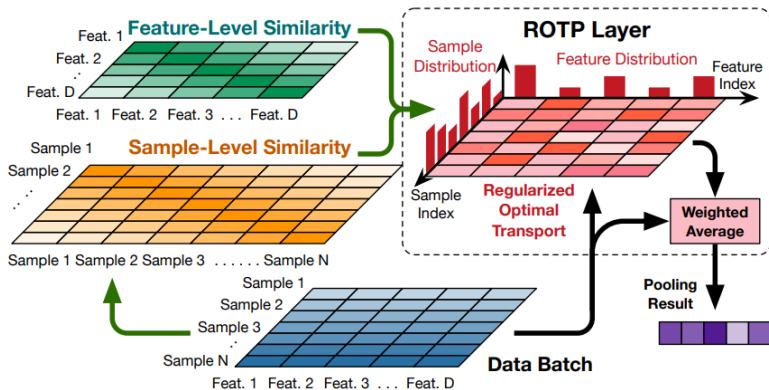
Design a global pooling layer = Determine the "sample-feature dimension" distribution \mathbf{P}

- ▶ Output energy maximization \Rightarrow OT term
- ▶ Consider correlation within data \Rightarrow GW structural regularizer
- ▶ Avoid sparse distribution \Rightarrow Bregman smoothness regularizer
- ▶ Prior of marginal distribution \Rightarrow unbalanced OT regularizer

$$\mathbf{P}_{\text{rot}}^*(\mathbf{X}; \theta) = \arg \min_{\mathbf{P} \in \Omega} \underbrace{\langle -\mathbf{X}, \mathbf{P} \rangle}_{\text{OT term}} + \underbrace{\alpha_0 \langle C(\mathbf{X}, \mathbf{P}), \mathbf{P} \rangle}_{\text{Structural Reg.}} + \underbrace{\alpha_1 R(\mathbf{P})}_{\text{Smoothness Reg.}} + \underbrace{\alpha_2 \text{KL}(\mathbf{P} \mathbf{1}_N | \mathbf{p}_0) + \alpha_3 \text{KL}(\mathbf{P}^T \mathbf{1}_D | \mathbf{q}_0)}_{\text{Marginal Reg.}}, \quad (2)$$

$$f_{\text{rot}}(\mathbf{X}; \theta) = (\mathbf{X} \odot \text{diag}^{-1}(\mathbf{P}_{\text{rot}}^*(\mathbf{X}; \theta) \mathbf{1}_N) \mathbf{P}_{\text{rot}}^*(\mathbf{X}; \theta)) \mathbf{1}_N. \quad (3)$$

Revisiting Global Pooling via Regularized OT



(a) Regularized optimal transport pooling (ROTP) layer

Revisiting Global Pooling via Regularized OT

Theoretical analysis

- ▶ Permutation-invariance
- ▶ Generalized framework for existing poolings

$f_{rot}(\mathbf{X}; \boldsymbol{\theta})$	<i>Mean-pooling</i>	<i>Max-pooling</i>	<i>Attention-pooling</i>
α_0	0	0	0
α_1	$\rightarrow \infty$	0	$\rightarrow \infty$
α_2	$\rightarrow \infty$	$\rightarrow \infty$	$\rightarrow \infty$
α_3	$\rightarrow \infty$	0	$\rightarrow \infty$
\mathbf{p}_0	$\frac{1}{D} \mathbf{1}_D$	$\frac{1}{D} \mathbf{1}_D$	$\frac{1}{D} \mathbf{1}_D$
\mathbf{q}_0	$\frac{1}{N} \mathbf{1}_N$	—	\mathbf{a}_X

Revisiting Global Pooling via Regularized OT

- ▶ Discrete Optimal Transport

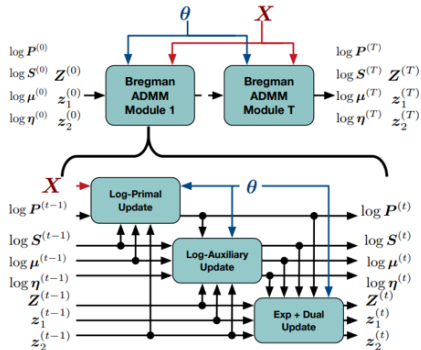
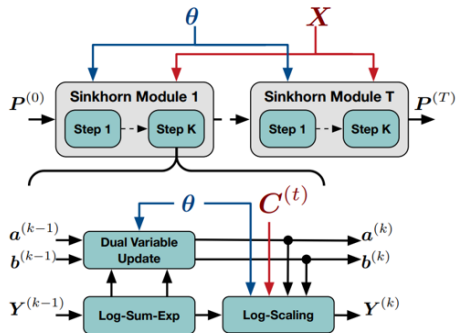
$$W_1(\boldsymbol{\mu}, \boldsymbol{\gamma}) := \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \langle \mathbf{D}, \mathbf{T} \rangle = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \sum_{v, v' \in \mathcal{V} \times \mathcal{V}} t_{vv'} d_{vv'}, \quad (4)$$

where $\Pi(\boldsymbol{\mu}, \boldsymbol{\gamma}) = \{\mathbf{T} \geq \mathbf{0} | \mathbf{T}\mathbf{1}_{|\mathcal{V}|} = \boldsymbol{\mu}, \mathbf{T}^\top \mathbf{1}_{|\mathcal{V}|} = \boldsymbol{\gamma}\}$

- ▶ Optimal transportation distances are parameterized distance
- ▶ Iterative computation methods: Sinkhorn, Baddmm

Revisiting Global Pooling via Regularized OT

Algorithmic modeling: Implement neural network layers by unrolling iterative optimization algorithmic.



Revisiting Global Pooling via Regularized OT

Experiments:

TABLE 1

Comparison on MIL accuracy \pm Std. (%) for different pooling layers.

Dataset	Messidor	Component	Function
D	687	200	200
#Positive bags	654	423	443
#Negative bags	546	2,707	4,799
#Instances	12,352	36,894	55,536
Min. bag size	8	1	1
Max. bag size	12	53	51
Add	74.33 \pm 2.56	93.35 \pm 0.98	96.26 \pm 0.48
Mean	74.42 \pm 2.47	93.32 \pm 0.99	96.28 \pm 0.66
Max	73.92 \pm 3.00	93.23 \pm 0.76	95.94 \pm 0.48
DeepSet	74.42 \pm 2.87	93.29 \pm 0.95	96.45 \pm 0.51
Mixed	73.42 \pm 2.29	93.45 \pm 0.61	96.41 \pm 0.53
GatedMixed	73.25 \pm 2.38	93.03 \pm 1.02	96.22 \pm 0.65
Set2Set	73.58 \pm 3.74	93.19 \pm 0.95	96.43 \pm 0.56
Attention	74.25 \pm 3.67	93.22 \pm 1.02	96.31 \pm 0.66
GatedAtt	73.67 \pm 2.23	93.42 \pm 0.91	96.51 \pm 0.77
DynamicP	73.16 \pm 2.12	93.26 \pm 1.30	96.47 \pm 0.58
GNP	73.54 \pm 3.68	92.86 \pm 1.96	96.10 \pm 1.03
OTK	74.78 \pm 2.89	93.19 \pm 0.93	96.31 \pm 1.02
SWE	74.46 \pm 3.72	93.32 \pm 1.26	96.42 \pm 0.88
ROTP _S	75.42 \pm 2.96	93.29 \pm 0.83	96.62 \pm 0.48
ROTP _{B-E} ($\alpha_0 = 0$)	74.83 \pm 2.07	93.16 \pm 1.02	96.17 \pm 0.43
ROTP _{B-Q} ($\alpha_0 = 0$)	75.08 \pm 2.06	93.13 \pm 0.94	96.09 \pm 0.46
ROTP _{B-E} (learn α_0)	75.33 \pm 1.96	93.16 \pm 1.08	96.22 \pm 0.44
ROTP _{B-Q} (learn α_0)	75.17 \pm 2.45	93.45 \pm 0.96	96.22 \pm 0.48

* The top-3 results are bolded and the best result is in red.

TABLE 2

Comparison on graph classification accuracy \pm Std. (%) for different pooling layers.

Dataset	NCII	PROTEINS	MUTAG	COLLAB	RD-T-B	RD-T-MSK	IMDB-B	IMDB-M
#Graphs	4,110	1,113	188	5,000	2,000	4,999	1,000	1,500
Average #Nodes	29.87	39.06	17.93	74.49	429.63	508.52	19.77	13.00
Average #Edges	32.30	72.82	19.79	2,457.78	497.75	594.87	96.53	65.94
#Classes	2	2	2	3	2	5	2	3
Add	67.96 \pm 0.43	72.97 \pm 0.54	89.05 \pm 0.86	71.06 \pm 0.43	80.00 \pm 1.49	50.16 \pm 0.97	70.18 \pm 0.87	47.56 \pm 0.56
Mean	64.82 \pm 0.52	66.09 \pm 0.64	86.53 \pm 1.62	72.35 \pm 0.44	83.62 \pm 1.18	52.44 \pm 1.24	70.34 \pm 0.38	48.65 \pm 0.91
Max	65.95 \pm 0.76	72.27 \pm 0.33	85.90 \pm 1.68	73.07 \pm 0.57	82.62 \pm 1.25	44.34 \pm 1.93	70.24 \pm 0.54	47.80 \pm 0.54
DeepSet	66.28 \pm 0.72	73.76 \pm 0.47	87.84 \pm 0.71	69.74 \pm 0.66	82.91 \pm 1.37	47.45 \pm 0.54	70.84 \pm 0.71	48.05 \pm 0.71
Mixed	66.46 \pm 0.74	72.25 \pm 0.45	87.30 \pm 0.87	73.22 \pm 0.35	84.36 \pm 2.62	46.67 \pm 1.63	71.28 \pm 0.26	48.07 \pm 0.25
GatedMixed	63.86 \pm 0.76	69.40 \pm 1.93	87.94 \pm 1.28	71.94 \pm 0.40	80.60 \pm 3.89	44.78 \pm 4.53	70.96 \pm 0.60	48.09 \pm 0.44
Set2Set	65.10 \pm 1.12	68.61 \pm 1.44	87.77 \pm 0.86	72.31 \pm 0.73	80.08 \pm 5.72	49.85 \pm 2.77	70.36 \pm 0.85	48.30 \pm 0.54
Attention	64.35 \pm 0.61	67.70 \pm 0.95	88.08 \pm 1.22	72.57 \pm 0.41	81.55 \pm 4.39	51.85 \pm 0.66	70.60 \pm 0.38	47.83 \pm 0.78
GatedAtt	64.66 \pm 0.52	68.16 \pm 0.90	86.91 \pm 1.79	72.31 \pm 0.37	82.55 \pm 1.96	51.47 \pm 0.82	70.52 \pm 0.31	48.67 \pm 0.35
DynamicP	62.11 \pm 0.77	65.86 \pm 0.85	85.40 \pm 2.81	70.78 \pm 0.88	67.51 \pm 1.82	32.11 \pm 3.85	69.84 \pm 0.73	47.59 \pm 0.48
GNP	68.20 \pm 0.48	73.44 \pm 0.61	88.37 \pm 1.25	72.80 \pm 0.58	81.93 \pm 2.23	51.80 \pm 0.61	70.34 \pm 0.83	48.85 \pm 0.81
ASAP	68.09 \pm 0.42	70.42 \pm 1.45	87.68 \pm 1.42	68.20 \pm 2.37	73.91 \pm 1.50	44.58 \pm 0.44	68.33 \pm 2.50	43.92 \pm 1.13
SAGP	67.48 \pm 0.65	72.63 \pm 0.44	87.88 \pm 2.22	70.19 \pm 0.55	74.12 \pm 2.86	46.00 \pm 1.74	70.34 \pm 0.74	47.04 \pm 1.22
OTK	67.96 \pm 0.55	69.09 \pm 0.76	86.90 \pm 1.83	71.35 \pm 0.91	74.28 \pm 1.39	50.57 \pm 1.20	70.94 \pm 0.79	48.41 \pm 0.89
SWE	68.06 \pm 0.98	70.52 \pm 1.22	85.68 \pm 2.07	72.17 \pm 1.29	79.30 \pm 3.94	51.11 \pm 1.53	70.34 \pm 1.05	48.93 \pm 1.34
WEGL	68.16 \pm 0.62	71.58 \pm 0.94	88.68 \pm 1.66	72.55 \pm 0.69	82.80 \pm 1.73	52.03 \pm 0.60	71.94 \pm 0.75	49.20 \pm 0.87
ROTP _S	68.27 \pm 1.06	73.10 \pm 0.22	88.84 \pm 1.21	71.20 \pm 0.55	81.54 \pm 1.38	51.00 \pm 0.61	70.74 \pm 0.80	47.87 \pm 0.43
ROTP _{B-E} ($\alpha_0 = 0$)	66.23 \pm 0.50	67.71 \pm 1.70	86.82 \pm 2.02	73.86 \pm 0.44	86.80 \pm 1.19	52.25 \pm 0.75	71.72 \pm 0.88	50.48 \pm 1.14
ROTP _{B-Q} ($\alpha_0 = 0$)	66.18 \pm 0.76	69.88 \pm 0.87	85.42 \pm 1.10	74.14 \pm 0.24	87.72 \pm 1.03	52.79 \pm 0.60	72.34 \pm 0.80	49.36 \pm 0.52
ROTP _{B-E} (learn α_0)	65.90 \pm 0.94	70.19 \pm 0.66	88.01 \pm 1.51	74.05 \pm 0.34	86.78 \pm 1.14	52.77 \pm 0.69	71.76 \pm 0.62	50.28 \pm 0.86
ROTP _{B-Q} (learn α_0)	65.96 \pm 0.32	70.12 \pm 1.17	86.79 \pm 1.81	74.27 \pm 0.47	88.67 \pm 0.99	52.84 \pm 0.60	71.78 \pm 1.00	49.44 \pm 0.46

* For each dataset, the top-3 results are bolded and the best result is in red.

TABLE 3

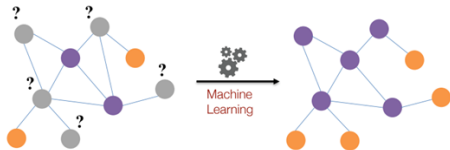
Comparisons on graph set classification accuracy \pm Std. (%) for different pooling layers.

Dataset	DECAGON DiBr-APND	DECAGON Anae-Fati	DECAGON PleuP-Diar	FEARS
#Graph sets	6,309	2,922	2,842	6,338
#Positive sets	3,189	1,526	1,422	3,169
Positive label	Difficulty breathing	Anaemia	Pleural pain	Non-mypath
#Negative sets	3,120	1,396	1,420	3,169
Negative label	Pressure decreased	Fatigue	Diarrhea	Myopathy
Set size	2	2	2	2~52
Add	50.86 \pm 0.97	63.15 \pm 1.79	62.32 \pm 1.08	75.89 \pm 1.33
Mean	51.10 \pm 1.09	61.95 \pm 2.60	61.30 \pm 2.68	72.42 \pm 1.51
Max	50.59 \pm 0.77	61.88 \pm 2.03	60.11 \pm 2.03	82.02 \pm 0.72
DeepSet	49.83 \pm 1.07	56.24 \pm 2.20	51.78 \pm 3.10	82.40 \pm 1.56
Mixed	51.13 \pm 0.99	63.83 \pm 1.19	60.91 \pm 2.12	81.54 \pm 1.13
GatedMixed	51.39 \pm 0.63	61.50 \pm 1.61	59.12 \pm 2.12	81.88 \pm 1.14
Set2Set	50.72 \pm 1.71	59.35 \pm 2.04	55.01 \pm 2.12	79.29 \pm 0.84
Attention	50.52 \pm 1.10	61.40 \pm 2.03	61.33 \pm 2.40	75.98 \pm 0.74
GatedAtt	50.74 \pm 0.61	62.15 \pm 0.77	58.80 \pm 1.18	75.84 \pm 1.29
DynamicP	51.01 \pm 1.88	55.93 \pm 1.56	52.58 \pm 2.91	74.00 \pm 1.61
GNP	50.00 \pm 1.88	53.98 \pm 4.34	52.58 \pm 4.68	62.71 \pm 15.55
ASAP	50.89 \pm 0.62	63.66 \pm 1.81	60.67 \pm 2.69	77.15 \pm 1.13
SAGP	49.87 \pm 0.77	63.62 \pm 1.28	59.86 \pm 2.43	77.29 \pm 1.04
OTK	50.96 \pm 1.11	63.68 \pm 1.59	61.66 \pm 2.39	79.40 \pm 1.08
SWE	51.05 \pm 2.15	63.21 \pm 2.02	61.37 \pm 3.13	80.64 \pm 1.86
WEGL	51.67 \pm 0.85	63.79 \pm 2.54	61.36 \pm 2.30	81.98 \pm 0.77
ROTP _S	51.96 \pm 0.71	62.91 \pm 1.13	59.40 \pm 0.90	79.75 \pm 0.71
ROTP _{B-E}	51.26 \pm 0.84	63.86 \pm 1.41	62.57 \pm 1.34	82.55 \pm 0.40
ROTP _{B-Q}	52.72 \pm 0.66	63.15 \pm 1.37	60.85 \pm 1.65	81.43 \pm 1.12

* The top-3 results are bolded and the best result is in red.

A Quasi-Wasserstein Loss for Learning Graph Neural Networks

Motivation: Eliminate Inconsistence of GNNs



$$\text{Loss function: } \min_{\theta} \sum_{v \in \mathcal{V}_L} \psi(g_v(\mathbf{X}, \mathbf{A}; \theta), \mathbf{y}_v). \quad (5)$$

- ▶ ψ is often implemented as cross-entropy loss, KL-divergence, Euclidean distance, and so on.
- ▶ The learning paradigm in (4) treats each node independently and evenly, inconsistent with the **non-i.i.d.** nature of graph-structured data.

A Quasi-Wasserstein Loss for Learning Graph Neural Networks

- ▶ Wasserstein distance:

$$QW(\hat{\mathbf{Y}}_{\mathcal{V}_L}, \mathbf{Y}_{\mathcal{V}_L}) = \sum_{c=1}^C W_1^{(P)}(\hat{\mathbf{y}}_{\mathcal{V}_L}^{(c)}, \mathbf{y}_{\mathcal{V}_L}^{(c)}) \quad (6)$$

- ▶ Wasserstein distance are parameterized distance

$$W_1(\boldsymbol{\mu}, \boldsymbol{\gamma}) := \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \langle \mathbf{D}, \mathbf{T} \rangle = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \sum_{v, v' \in \mathcal{V} \times \mathcal{V}} t_{vv'} d_{vv'}, \quad (7)$$

where $\Pi(\boldsymbol{\mu}, \boldsymbol{\gamma}) = \{\mathbf{T} \geq \mathbf{0} \mid \mathbf{T}\mathbf{1}_{|\mathcal{V}|} = \boldsymbol{\mu}, \mathbf{T}^\top \mathbf{1}_{|\mathcal{V}|} = \boldsymbol{\gamma}\}$

- ▶ Optimal Transport on Graphs– minimum-cost flow problem^[1]

$$W_1(\boldsymbol{\mu}, \boldsymbol{\gamma}) = \min_{\mathbf{f} \in \Omega_{|\mathcal{E}|}} \|\text{diag}(\mathbf{w})\mathbf{f}\|_1 \quad \text{s.t. } \mathbf{S}\mathbf{v}\mathbf{f} = \boldsymbol{\gamma} - \boldsymbol{\mu}, \quad (8)$$

where

$$s_{ve} = \begin{cases} 1 & \text{if } v \text{ is “head” of edge } e \\ -1 & \text{if } v \text{ is “tail” of edge } e \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

Montacer Essid and Justin Solomon. Quadratically regularized optimal transport on graphs.391 SIAM Journal on Scientific Computing, 40(4):A1961–A1986, 2018.

Enrico Facca and Michele Benzi. Fast iterative solution of the optimal transport problem on393 graphs. SIAM Journal on Scientific Computing, 43(3):A2295–A2319, 2021

A Quasi-Wasserstein Loss for Learning Graph Neural Networks

- Minimum-cost flow problem:

$$\begin{aligned}
 QW(\hat{\mathbf{Y}}_{\mathcal{V}_L}, \mathbf{Y}_{\mathcal{V}_L}) &= \sum_{c=1}^C W_1^{(P)}(\hat{\mathbf{y}}_{\mathcal{V}_L}^{(c)}, \mathbf{y}_{\mathcal{V}_L}^{(c)}) \\
 &= \sum_{c=1}^C \min_{\mathbf{f}^{(c)} \in \Omega} \|\text{diag}(\mathbf{w})\mathbf{f}^{(c)}\|_1 \quad \text{s.t. } \mathbf{S}_{\mathcal{V}_L}\mathbf{f}^{(c)} = \mathbf{y}_{\mathcal{V}_L}^{(c)} - \hat{\mathbf{y}}_{\mathcal{V}_L}^{(c)} \\
 &= \min_{\mathbf{F} \in \Omega_C} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1 \quad \text{s.t. } \mathbf{S}_{\mathcal{V}_L}\mathbf{F} = \mathbf{Y}_{\mathcal{V}_L} - g_{\mathcal{V}_L}(\mathbf{A}, \mathbf{X}; \theta),
 \end{aligned} \tag{10}$$

- Learning Paradigm:

$$\min_{\theta} \underbrace{\min_{\mathbf{F} \in \Omega} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1 \quad \text{s.t. } \mathbf{S}_{\mathcal{V}_L}\mathbf{F} = \mathbf{Y}_{\mathcal{V}_L} - g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta)}_{QW(\hat{\mathbf{Y}}_{\mathcal{V}_L}, \mathbf{Y}_{\mathcal{V}_L})}. \tag{11}$$

\Downarrow

$$\min_{\theta, \mathbf{F} \in \Omega} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1 + \lambda \underbrace{B_{\phi}(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_{\mathcal{V}_L}\mathbf{F}, \mathbf{Y}_{\mathcal{V}_L})}_{\sum_{v \in \mathcal{V}_L} \psi(g_v(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_v\mathbf{F}, \mathbf{y}_v)}. \tag{12}$$

- A New Transductive Prediction Paradigm:

$$\tilde{\mathbf{y}}_v := g_v(\mathbf{X}, \mathbf{A}(\mathbf{F}^*; \xi^*); \theta^*) + \mathbf{S}_v\mathbf{F}^*, \tag{13}$$

A Quasi-Wasserstein Loss for Learning Graph Neural Networks

Experiments:

Table 2: Comparisons on node classification accuracy (%) on homophilic graphs.

Model	Method	Cora	Citeseer	Pubmed	Computers	Photo	Improve
GCN	(\mathbb{I})	87.14 \pm 1.01	79.86 \pm 0.67	86.74 \pm 0.27	83.32 \pm 0.33	88.26 \pm 0.73	
	(\mathbb{I})+LPA	86.34 \pm 1.45	78.51 \pm 1.22	84.72 \pm 0.70	82.48 \pm 0.69	88.10 \pm 1.31	-1.07
	QW	86.95 \pm 1.12	81.30 \pm 0.37	87.89 \pm 0.44	88.39 \pm 0.55	93.80 \pm 0.37	+2.60
APNP	(\mathbb{I})	88.14 \pm 0.73	80.47 \pm 0.74	88.12 \pm 0.31	85.32 \pm 0.37	88.51 \pm 0.31	
	QW	88.65 \pm 1.00	80.94 \pm 0.61	89.39 \pm 0.31	84.69 \pm 0.46	93.25 \pm 0.38	+0.94
BernNet	(\mathbb{I})	88.28 \pm 1.00	79.81 \pm 0.79	88.87 \pm 0.38	87.61 \pm 0.46	93.68 \pm 0.28	
	QW	89.03\pm0.76	81.35\pm0.71	89.03 \pm 0.38	89.14 \pm 0.39	94.28 \pm 0.44	+0.92
ChebNetII	(\mathbb{I})	88.26 \pm 0.89	80.00 \pm 0.74	88.57 \pm 0.36	86.58 \pm 0.71	93.50 \pm 0.34	
	QW	87.98 \pm 0.80	79.47 \pm 0.70	89.15\pm0.44	89.52\pm0.54	94.84\pm0.37	+0.81

Table 3: Comparisons on node classification accuracy (%) on heterophilic graphs.

Model	Method	Squirrel	Chameleon	Actor	Texas	Cornell	Improve
GCN	(\mathbb{I})	46.55 \pm 1.15	63.57 \pm 1.16	34.00 \pm 1.28	77.21 \pm 3.28	61.91 \pm 5.11	
	(\mathbb{I})+LPA	44.81 \pm 1.81	60.90 \pm 1.63	32.43 \pm 1.59	78.69 \pm 6.47	68.72 \pm 5.95	+0.46
	QW	51.04 \pm 0.51	67.77 \pm 0.92	38.09 \pm 0.50	84.10 \pm 2.95	84.26 \pm 2.98	+8.40
APNP	(\mathbb{I})	36.15 \pm 0.75	52.93 \pm 1.71	40.46 \pm 0.64	91.31 \pm 1.97	87.66 \pm 2.13	
	QW	37.11 \pm 0.60	53.76 \pm 1.25	40.78 \pm 0.74	91.48 \pm 2.30	87.87 \pm 2.34	+0.50
BernNet	(\mathbb{I})	51.15 \pm 1.09	67.96 \pm 1.05	40.72 \pm 0.80	93.28 \pm 1.48	90.21 \pm 2.35	
	QW	53.29 \pm 0.65	70.96 \pm 1.31	40.91 \pm 0.71	93.44 \pm 1.80	90.85\pm2.34	+1.23
ChebNetII	(\mathbb{I})	57.78 \pm 0.84	71.71 \pm 1.40	40.70 \pm 0.77	92.79 \pm 1.48	88.94 \pm 2.78	
	QW	59.55\pm0.86	74.05\pm1.12	41.37\pm0.67	93.93\pm0.98	87.23 \pm 3.62	+0.84

Table 8: Comparisons on node classification accuracy (%) on homophilic graphs.

Model	Method	Cora	Citeseer	Pubmed	Computers	Photo	Improve
MLP	(\mathbb{I})	77.16 \pm 1.10	76.71 \pm 0.86	86.14 \pm 0.32	84.32 \pm 0.53	89.42 \pm 0.39	
	QW	76.77 \pm 0.99	77.74 \pm 0.60	86.56 \pm 0.41	82.39 \pm 0.40	89.51 \pm 0.34	-0.216
GAT	(\mathbb{I})	89.20 \pm 0.79	80.75 \pm 0.78	87.42 \pm 0.33	90.08 \pm 0.36	94.38 \pm 0.25	
	QW	88.56 \pm 0.92	80.19 \pm 0.64	88.38 \pm 0.23	90.41 \pm 0.28	94.65 \pm 0.24	+0.05

Table 9: Comparisons on node classification accuracy (%) on heterophilic graphs.

Model	Method	Squirrel	Chameleon	Actor	Texas	Cornell	Improve
MLP	(\mathbb{I})	32.19 \pm 0.77	48.21 \pm 1.47	40.61 \pm 0.60	91.80 \pm 1.31	88.30 \pm 2.55	
	QW	34.41 \pm 0.47	49.61 \pm 1.31	40.63 \pm 0.51	91.15 \pm 2.30	88.72 \pm 2.77	+0.68
GAT	(\mathbb{I})	48.20 \pm 1.67	64.31 \pm 2.01	35.68 \pm 0.60	80.00 \pm 3.11	68.09 \pm 2.13	
	QW	55.03 \pm 1.35	67.35 \pm 1.42	33.83 \pm 1.07	80.33 \pm 1.97	70.21 \pm 2.13	+2.09

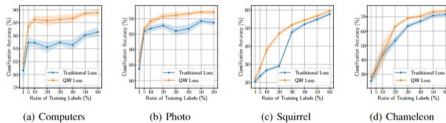


Figure 2: Illustrations of the learning methods' performance given different amounts of labeled nodes.