



AugWard: Augmentation-Aware Representation Learning for Accurate Graph Classification

Minjun Kim¹, Jaehyeon Choi¹, SeungJoo Lee¹, Jinhong Jung^{2*}, and U Kang^{1*}

¹Seoul National University, ²Soongsil University *Correspond to: ukang@snu.ac.kr, jinhong@ssu.ac.kr







- We propose AugWard, an augmentation-aware learning framework for accurate graph classification
 - AugWard enriches graph representations by capturing "augmentation-induced differences"
- AugWard is easily integrated with any method, enhancing their accuracy across various settings
 Supervised, semi-supervised, transfer learning







Introduction

- Preliminaries
- Proposed Method
- Experiments
- Conclusion





Graph Classification

- Task: Classify a graph into pre-defined classes based on its structural properties and features
- Graph Neural Networks (GNNs) capture higherorder structures for accurate classification



[1] M. Do et al., (2023) "Two-Stage Training of Graph Neural Networks for Graph Classification" Neural Process Letters 55, pp. 2799–2823 Minjun Kim (SNU)





Problem Definition Graph Classification

Given

- A set $G = \{G_1, \dots, G_N\}$ of N distinct graphs
- A set $Y = \{y_1, \dots, y_N\}$ of corresponding labels
- A set C of classes

Predict

• $\forall G_i \in G, y \in C$, the probability $P(y|G_i)$





Graph Augmentation

- Graph augmentation generates variants of original graphs while preserving their labels
 - Mitigate the common issue of overfitting



Model-agnostic graph augmentation





Rich decision boundary learned by f

[1] J. Yoo et al., (2022) "Model-Agnostic Augmentation for Accurate Graph Classification" WWW 2022

Minjun Kim (SNU)



- Existing methods are suboptimal due to 2 major limitations by simplistic adaptation of augmentation:
 - Limitation 1. Ignorance of difference between original and augmented graphs
 - Limitation 2. Deceptive assumption that the perturbation ratio p ensures similarity among augmented graphs









- Introduction
- Preliminaries
 - Proposed Method
 - Experiments
 - Conclusion





• Graph Neural Networks (GNNs) jointly train the encoder f_{θ} and the classifier g_{ϕ} for classification:

$$\mathbf{z}_{\mathcal{G}} = f_{\theta}(\mathcal{G}), \qquad \mathbf{p}_{\mathcal{G}} = g_{\phi}(\mathbf{z}_{\mathcal{G}})$$

- $\mathbf{z}_{\mathcal{G}} \in \mathbb{R}^{d}$: representation of a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X})$
- $\mathbf{p}_{\mathcal{G}} \in \mathbb{R}^{|\mathcal{C}|}$: predicted probabilities (*i*-th entry of $\mathbf{p}_{\mathcal{G}}$) = $P(y = i|\mathcal{G})$

$$\mathbf{z}_{\mathcal{G}} = \operatorname{READOUT}\left(\left\{\mathbf{h}_{u}^{(l)} | u \in \mathcal{V}, l \in [L]\right\}\right),\\ \tilde{\mathbf{h}}_{u}^{(l)} = \operatorname{AGGREGATE}\left(\left\{\mathbf{h}_{v}^{(l-1)} : v \in \mathcal{N}_{u}\right\}\right), \mathbf{h}_{u}^{(l)} = \operatorname{COMBINE}\left(\mathbf{h}_{u}^{(l-1)}, \tilde{\mathbf{h}}_{u}^{(l)}\right).$$

• $\mathbf{h}_{u}^{(0)} = \mathbf{X}_{u}, \ [L] = \{0, \dots, L\}, \ \mathcal{N}_{u}: \text{ set of neighbors for node } u$





Graph Augmentation

• Augmented graph \mathcal{G}^+ is randomly sampled as:

$$\mathcal{G}^+ \sim \mathcal{T}_p(\mathcal{G}^+|\mathcal{G})$$

- *G*: original graph
- $\mathcal{T}_p(\cdot | \mathcal{G})$: augmentation distribution conditioned on \mathcal{G}
- p: perturbation ratio (amount of change from G)
- Various designs for \mathcal{T}_p
 - 1. Drop-based methods: remove or mask attributes of nodes and edges according to ratio p
 - 2. Mixup-based methods: fuse two graphs by ratio p







Fused Gromov-Wasserstein Distance

FGWD measures the distance between two graphs G = (V, E, X) and G⁺ = (V, E, X⁺) by combining feature-level and structure-level differences

$$\operatorname{FGWD}_{\alpha}(\mathcal{G},\mathcal{G}^+) = \min_{\pi \in \Pi(\mu,\nu)} \alpha \cdot \frac{WD(\mathbf{X},\mathbf{X}^+,\pi)}{WD(\mathbf{X},\mathbf{X}^+,\pi)} + (1-\alpha) \ GWD(\mathcal{E},\mathcal{E}^+,\pi)$$

Wasserstein Distance (Feature-level distance) Gromov-Wasserstein Distance (Structure-level distance)

GWD

WD

- α: balancing hyperparameter
- π : coupling matrix
- $\Pi(\mu, \nu)$: a set of all possible matchings of distributions $\mu \in \mathbb{R}^{|\mathcal{V}|}$ and $\nu \in \mathbb{R}^{|\mathcal{V}^+|}$

[1] T. Vayer et al., (2020). "Fused Gromov-Wasserstein distance for structured objects". Algorithms, 13(9), 212.

Minjun Kim (SNU)







- Introduction
- Preliminaries
- Proposed Method
 - Experiments
 - Conclusion







AugWard for accurate graph classification

- Considers the "augmentation-induced difference"
- Applicable to any augmentation \mathcal{T}_p / encoder f_{θ} / classifier g_{ϕ}



Minjun Kim (SNU)







AugWard for accurate graph classification

- Idea 1. Augmentation-aware training
- Idea 2. Graph distance-based difference
- Idea 3. Consistency regularization







The impact of proposed ideas

- Ideas 1 and 2 enrich the representation of each graph
- Idea 3 regulates the classifier for consistent predictions





Loss Function



- Combining all 3 ideas, we get the objective function
 - AugWard supports various learning paradigms with baseline loss L_{base}:

Baseline loss (e.g., cross-entropy)

$$\mathcal{L}(\mathcal{G}, \mathcal{G}^+, y) = \mathcal{L}_{\text{base}}(\mathcal{G}, \mathcal{G}^+, y) + \mathcal{L}_{\text{AugWard}}(\mathcal{G}, \mathcal{G}^+),$$

$$\mathcal{L}_{\text{AugWard}}(\mathcal{G}, \mathcal{G}^+) = \lambda_{\text{aware}} \mathcal{L}_{\text{aware}}(\mathcal{G}, \mathcal{G}^+) + \lambda_{\text{cr}} \mathcal{L}_{\text{cr}}(\mathcal{G}, \mathcal{G}^+)$$

Idea 1. Augmentation-aware training Idea 2. Graph distance-based difference Idea 3. Consistency Regularization





17

Main Ideas

1. Augmentation-aware training

Challenge: Capturing augmentation-induced difference

- Existing methods ignore the difference between original and augmented graphs, or augmentation-induced difference
- Pearson Correlation Coefficient (PCC) close to 0 indicates no strong correlation Perturbation Ratio







Main Ideas

1. Augmentation-aware training

Idea: Augmentation-aware training

- Intuition. Encourage the encoder f_{θ} to align representation-level and graph-level differences
- Train a neural network h_{ω} by optimizing \mathcal{L}_{aware} :

$$\mathcal{L}_{\text{aware}}(\mathcal{G},\mathcal{G}^+) = \left\| h_{\omega} (\mathbf{z}_{\mathcal{G}}, \mathbf{z}_{\mathcal{G}^+}) - \mathcal{D}(\mathcal{G}, \mathcal{G}^+) \right\|^2$$

• $\mathcal{D}(\cdot,\cdot)$: graph-level distance

"Which metric is suitable for $\mathcal{D}(\mathcal{G},\mathcal{G}^+)$?"

** A fully connected layer h_{ω} with the concatenation of $\mathbf{z}_{\mathcal{G}}$ and $\mathbf{z}_{\mathcal{G}^+}$ as input









2. Graph Distance-based Difference

- Challenge: Measuring the difference gained from graph augmentation
 - Generated graphs at a fixed perturbation ratio p exhibit significant variations due to inherent randomness



** Euclidean $\|\mathbf{z}_{\mathcal{G}} - \mathbf{z}_{\mathcal{G}^+}\|_2^2$ distances between \mathcal{G} and 100 augmented graphs \mathcal{G}^+ with p fixed at 0.2 and 0.4 Minjun Kim (SNU) 19







2. Graph Distance-based Difference

- Idea: Graph Distance-based Difference
 - Fused Gromov-Wasserstein Distance (FGWD): optimizes both differences in structures and features

$$\mathcal{L}_{\text{aware}} = \left\| h_{\omega} (\mathbf{z}_{\mathcal{G}}, \mathbf{z}_{\mathcal{G}^+}) - \text{FGWD}_{\alpha} (\mathcal{G}, \mathcal{G}^+) \right\|^2$$

Difference in either structure or feature leads to a significant distinction in their chemical type



** (a) and (b) share graph features, while (b) and (c) exhibit identical graph structures Minjun Kim (SNU)





Main Ideas

3. Consistency Regularization

Motivation: Training robust classifier

- Given distinguishable representations, training the classifier robustly is crucial for better generalization
- In node classification^[1], matching predictions from different representations (same label) improves the generalization

Idea: Consistency regularization

Matching two predictions $p_{\mathcal{G}}$ and $p_{\mathcal{G}^+}$

$$\mathcal{L}_{cr} = H(\mathbf{p}_{\mathcal{G}}, \mathbf{p}_{\mathcal{G}^+}) = -\sum_{i=1}^{|\mathcal{C}|} P(y = i|\mathcal{G}) \cdot \log P(y = i|\mathcal{G}^+)$$

•
$$H(\cdot,\cdot)$$
: cross-entropy loss

[1] W. Feng et al., (2020) "Graph Random Neural Network for Semi-Supervised Learning on Graphs" NeurIPS 2020







- Introduction
- Preliminaries
- Proposed Method
- Experiments
 - Conclusion







Datasets

- Supervised, semi-supervised learning: 10 TUDatasets
- Transfer learning: ZINC15 → 8 MoleculeNet datasets

Augmentations

- Drop-based: NodeDrop, EdgeDrop, AttrMask, Subgraph, GraphAug
- Mixup-based: SubMix, S-Mixup

Baselines

- Model: a 4-layered GIN
- Semi-supervised (10%): Infograph, GraphCL, CuCo, GCL-SPAN
- Transfer learning: ContextPred, GraphCL, MGSSL, GraphMAE





- We perform experiments on the following questions:
 - Q1. Accuracy in supervised graph classification
 - Q2. Accuracy in semi-supervised graph classification
 - Q3. Representation transferability
 - Q4. Runtime analysis
 - Q5. Ablation study







- 1. Accuracy in supervised graph classification
- AugWard consistently enhances classification accuracy
 - Up to 2.13%p in average accuracy









2. Accuracy in semi-supervised graph classification

- AugWard is also beneficial in semi-supervised setting
 - Up to 1.52%p increase in average accuracy







Experiments

3. Representation Transferability

- AugWard offers more expressive representations
 - Improves the performance of transfer learning models; up to 3.71%p in average accuracy







Experiments 4. Runtime Analysis

- The overhead from AugWard is marginal
 - Computing FGWD takes 4.89% of total time in average







Experiments 5. Ablation Study

- All ideas contribute to the enhanced performance
 - Considering the *augmentation-induced difference* is beneficial, even with simple heuristics

	Variants	Ideas	Avg.	Imp.
	A: $GIN + NodeDrop$	Existing	61.58	-
Naïve metrics for graph-level distance	$\mathrm{A}+p$	I1	61.90	+0.32
	A + NFs	I1	62.50	+0.92
	A + AMs	I1	62.55	+0.98
	A + Edge Jaccard	I1	62.53	+0.95
FGWD (Proposed)	A + FGWD	I1+I2	63.04	+1.46
	$A + \mathbf{AugWard}$	I1+I2+I3	63.58	+2.00







- Introduction
- Preliminaries
- Proposed Method
- Experiments
- Conclusion







- We propose AugWard for graph classification
 - AugWard considers the augmentation-induced differences
- Main ideas
 - Ideas 1 & 2. Augmentation-aware training with FGWD
 - Ideas 3. Consistency regularization

Experiments

 AugWard consistently enhances classification performance across various learning settings





Thank you !

Minjun Kim (minjun.kim@snu.ac.kr)





