

snmsung Research



Datasets, Tasks, and Training Methods for Large-scale Hypergraph Learning

> Sunwoo Kim*, Dongjin Lee*, Yul Kim, Jungho Park, Taeho Hwang, Kijung Shin

Overview

- Proposed Tasks
- Proposed Datasets
- Proposed Training Method
- O Experimental Results

• Conclusions



Group Interactions are EVERYWHERE!





Protein Interaction



 Excited for our new monarch!
 ~

 To
 (mccarth@gmail.comx)
 Cc

 Bcc
 (jaffe@icloud.comx)
 (mhoffman@optonline.net x)
 Cc

 Bcc
 (jaffe@icloud.comx)
 (mhoffman@optonline.net x)
 Cc

 Excited for our new monarch!
 Excited for our new monarch!
 Excited for our new monarch!

Email

Sunwoo Kim

Datasets, Tasks, and Training Methods for Large-scale Hypergraph Learning

Group Interactions → HYPERGRAPH!

- Group interactions can be modeled as a hypergraph.
- A hypergraph consists of a node set and a hyperedge set.



Applications of Hypergraphs

• Hypergraphs are widely used in various domains.



Challenges in Hypergraph Learning

• The current hypergraph representation learning field has three challenges:



Limited downstream task

Small benchmark datasets

No scalable training strategy

Our Contribution

• Our contribution is three-fold.



Overview

- O Proposed Tasks
- Proposed Datasets
- Proposed Training Method
- O Experimental Results

Conclusions



Proposed Tasks

- We propose two pair-level downstream tasks in hypergraphs.
 - Task 1) Hyperedge disambiguation, classifying a pair of hyperedges.
 - Task 2) Local clustering, classifying a pair of nodes.



Proposed Task 1: Formulation

- Setting: Certain hyperedges are split into two (or more) hyperedges.
- Goal: Identifying whether a pair of hyperedges are originally identical hyperedges or not.



Proposed Task 1: Application

- One application of task 1 is author disambiguation.
 - Goal: Identifying whether two authors are identical authors or not.
 - Hypergraph formulation: Nodes (publications) & Hyperedges (authors)



Proposed Task 2: Formulation

- Setting: Each node belongs to at least one cluster.
- Goal: Identifying whether a pair of nodes belong to the same cluster or not.

Proposed Task 2: Application

- One application of task 2 is device-household matching.
 - · Goal: Identifying whether two devices belong to the same household or not.
 - Hypergraph formulation: Nodes (devices) & Hyperedges (IP access) & Clusters (households).
 Access to Accessed devices

Sunwoo Kim

Datasets, Tasks, and Training Methods for Large-scale Hypergraph Learning

Overview

- Proposed Tasks
- **O** Proposed Datasets
- Proposed Training Method
- O Experimental Results

Conclusions

Proposed Dataset: Overview

• We propose two large-scale co-authorship hypergraph datasets:

Dataset	# of Nodes	# of Hyperedges	# Features	# of Classes
AMiner	13,262,573	22,552,647	300	257
MAG	27,320,375	30,175,013	300	247

Proposed Dataset: Details

- Source of datasets: Open Academic Graph 2.1
- Nodes: Publications
- Hyperedges: Authors

Proposed Dataset: Details

- Node features: An embedding of corpus, which consists of title and keywords, transformed via Sentence2Vec (S2V).
- Node labels: Research field of the corresponding publication.

<u>Title:</u> Temporal patterns of real-world group interactions.

Keywords: Hypergraph, Group

interaction, Time series

Field: Data Mining & Analysis

Overview

- Proposed Tasks
- Proposed Datasets
- O Proposed Training Method
- O Experimental Results

Conclusions

Challenges in Training HNNs on Large-Scale HG

- <u>Challenge 1</u>) Considerable time and space cost.
 - Full-batch training: Can't load the entire hypergraph into GPU memory.
 - Mini-batch training: Computational overhead occurs in obtaining subhypergraphs and loading them into GPU memory.

Proposed Learning Method: Solution 1

- <u>Solution 1</u>) Partitioning the input hypergraph.
 - A partitioning strategy is efficient in training graph neural networks (Chiang et al. 2019).
 - We use partitioning, treating each hypergraph partition as a single mini-batch.

Challenges in Training HNNs on Large-Scale HG

- <u>Challenge 2</u>) Loss of pair-label supervision due to partitioning.
 - When we process a single partition at a time, certain pair-labels cannot be utilized.

Proposed Learning Method: Solution 2

- <u>Solution 2</u>) Utilizing <u>Contrastive Learning (CL)</u>.
 - CL can make effective representations under label-scarce scenarios (Chen et al. ICML 2020, You et al. NeurIPS 2020).

Proposed Learning Method: Solution 2

- Overview of CL is as follows:
 - CL first creates multiple views of a hypergraph.
 - Then, it lets HNN learn the agreement of two views.

- We propose PCL: Partitioning-based Contrastive Learning.
 - We first divide the input hypergraph into several partitions.
 - Then, we train HNN by CL, by regarding each partition as a minibatch of CL.

• Specifically, we do CL for each partition, for the given number of epochs.

- When the CL process is done, we obtain node embeddings with HNN.
- Here, we do not additionally train HNN to:
 - Reduce cost of message passing process.
 - Utilize pair of nodes (hyperedges) in different partitions.

• With node embeddings and a classifier, we perform pair-level tasks.

- We propose two additional tools for PCL.
 - P-IOS: Partitioning technique that mitigates topological information loss.
 - PINS: CL technique that encourages HNN to learn inter-partition dissimilarity.

- <u>P-IOS</u>: <u>Partitioning with the Inclusion of Outsider Set.</u>
 - Certain hyperedges are destroyed during the partitioning process.

Datasets, Tasks, and Training Methods for Large-scale Hypergraph Learning

- <u>P-IOS</u>: <u>Partitioning with the Inclusion of Outsider Set.</u>
 - Certain hyperedges are destroyed during the partitioning process.
 - P-IOS recovers destroyed hyperedges and preserves every 1-hop neighboring node.

Datasets, Tasks, and Training Methods for Large-scale Hypergraph Learning

- <u>PINS</u>: <u>Previous part</u>tion's <u>Negative</u> <u>Samples</u>.
 - Basic PCL utilizes a single partition at a time, and HNN cannot learn

dissimilarity between nodes at different partitions.

- <u>PINS:</u> <u>P</u>revious part<u>I</u>tion's <u>N</u>egative <u>S</u>amples.
 - Basic PCL utilizes a single partition at a time, and HNN can't learn dissimilarity between nodes at different partitions.
 - PINS uses embeddings of previous partitions as negative samples.

Overview

- Proposed Tasks
- Proposed Datasets
- Proposed Training Method
- O Experimental Results
- O Conclusions

Research Questions

- We aim to analyze the following four points:
 - <u>RQ1</u>) Accuracy of PCL for the proposed pair-level tasks.
- <u>RQ2</u>) Effectiveness of PINS.
- <u>RQ3</u>) Memory consumption of PINS.

• <u>RQ4</u>) Effectiveness of P-IOS.

Used Datasets

• We utilize 5 real-world hypergraph datasets.

	Dataset	# Nodes	# Hyperedges	# Features	# Classes
bc	DBLP	41,302	22,263	1,425	6
kistiı	Trivago	172,738	233,202	300	160
ш Т	OGBN-MAG	736,389	1,134,649	128	349
osec	AMiner	13,262,573	22,552,647	300	257
Prop	MAG	27,320,375	30,175,013	300	257

Sunwoo Kim

Microsoft

Research

• RQ1) Accuracy of PCL on Task 1 (Hyperedge disambiguation).

				-								
		DB	BLP	Trivago		OGBN	OGBN-MAG		AMiner		AG	Avg
Data Type	Methods	AP	AUROC	AP	AUROC	AP	AUROC	AP	AUROC	AP	AUROC	Rank
Only X	MLP	65.2 ± 3.4	62.2 ± 3.4	60.2 ± 1.6	60.7 ± 2.5	73.1 ± 2.6	71.9 ± 2.8	91.7 ± 0.3	$\underline{92.9\pm0.3}$	$\underline{91.2\pm0.6}$	$\underline{92.6\pm0.6}$	9.8
	GCN	83.6 ± 1.2	84.6 ± 1.5	61.6 ± 0.8	58.1 ± 2.4	80.1 ± 1.7	79.8 ± 2.0	OOM	OOM	OOM	OOM	9.7
Full	GAT	83.5 ± 1.0	84.6 ± 1.0	61.6 ± 0.5	60.0 ± 2.7	76.5 ± 2.3	75.7 ± 3.0	OOM	OOM	OOM	OOM	10.3
Graph	BGRL	84.5 ± 1.2	84.3 ± 1.7	72.0 ± 1.8	72.2 ± 2.2	83.1 ± 0.6	83.9 ± 0.5	OOM	OOM	OOM	OOM	8
	GGD	80.9 ± 1.6	80.6 ± 2.1	71.8 ± 0.9	72.5 ± 1.6	76.1 ± 1.1	76.1 ± 1.3	OOM	OOM	OOM	OOM	9.7
	HGNN	77.3 ± 2.5	80.0 ± 3.1	72.1 ± 5.1	76.6 ± 1.4	91.2 ± 0.7	92.9 ± 0.6	OOM	OOM	OOM	OOM	7.7
Full	UniGCNII	78.1 ± 3.4	77.1 ± 3.9	71.8 ± 1.5	75.6 ± 2.3	86.9 ± 5.6	88.3 ± 6.7	OOM	OOM	OOM	OOM	8.7
Hypergraph	ALLSET	53.1 ± 1.2	55.8 ± 2.2	53.7 ± 1.2	56.8 ± 2.2	65.3 ± 1.4	73.4 ± 1.6	OOM	OOM	OOM	OOM	13.1
	HCL	$\textbf{91.0} \pm \textbf{1.7}$	$\textbf{92.6} \pm \textbf{2.1}$	73.7 ± 0.7	78.8 ± 1.1	$\textbf{94.2}\pm\textbf{0.4}$	$\textbf{95.5} \pm \textbf{0.3}$	OOM	OOM	OOM	OOM	4.5
άτ.	GCN	84.5 ± 1.6	85.9 ± 1.6	65.9 ± 1.9	66.7 ± 2.0	82.3 ± 2.6	83.4 ± 2.8	81.2 ± 0.6	81.5 ± 0.6	OOM	OOM	7.9
Partitioned	GAT	84.7 ± 1.3	86.4 ± 1.2	62.2 ± 3.1	62.5 ± 4.6	78.0 ± 1.8	79.6 ± 1.8	81.3 ± 0.5	81.0 ± 0.6	OOM	OOM	8.2
Graph	BGRL	85.8 ± 1.7	85.6 ± 1.7	$\underline{82.4\pm2.2}$	$\underline{83.5\pm2.1}$	91.4 ± 0.5	92.8 ± 0.5	90.1 ± 0.8	90.9 ± 0.8	89.7 ± 0.7	90.8 ± 0.9	3.9
	GGD	80.1 ± 2.4	80.8 ± 2.3	76.9 ± 1.6	78.3 ± 1.5	88.7 ± 1.0	90.2 ± 0.9	82.5 ± 1.2	83.6 ± 1.5	80.1 ± 1.3	81.2 ± 1.4	6.2
Dantitionad	HGNN	84.1 ± 2.0	86.2 ± 1.8	71.8 ± 0.8	75.8 ± 1.0	85.9 ± 1.1	88.1 ± 1.0	91.3 ± 0.3	92.5 ± 0.3	89.2 ± 0.7	91.6 ± 0.5	5.4
Partitioned	UniGCNII	79.9 ± 1.8	80.0 ± 1.8	71.5 ± 2.6	74.2 ± 3.5	77.4 ± 0.9	75.7 ± 1.2	91.8 ± 0.5	92.1 ± 0.5	90.4 ± 0.4	91.9 ± 0.1	7.9
nypergraph	ALLSET	54.2 ± 1.6	57.6 ± 2.7	53.6 ± 1.2	56.7 ± 2.0	59.3 ± 1.6	65.6 ± 2.4	61.5 ± 1.8	68.6 ± 2.5	OOM	OOM	13.2
Partitioned Hypergraph	Basic PCL (proposed)	87.1 ± 1.9	$\underline{88.3\pm2.0}$	$\textbf{88.2}\pm\textbf{0.9}$	$\textbf{88.9} \pm \textbf{0.7}$	94.1 ± 0.4	$\underline{95.1\pm0.3}$	95.5 ± 0.7	$\textbf{96.0}\pm\textbf{0.8}$	$\textbf{96.2}\pm\textbf{0.3}$	$\textbf{96.9}\pm\textbf{0.3}$	1.4
												(

*bold: best

**<u>underline: second-best</u>

• RQ1) Accuracy of PCL on Task 2 (Local clustering).

Data Type	Methods	DB AP	LP AUROC	Triv AP	^z ago AUROC	OGBN AP	-MAG AUROC	AM AP	iner AUROC	MA AP	AG AUROC	Avg. Rank
Only X	MLP	52.2 ± 2.3	51.3 ± 3.5	50.6 ± 0.4	50.9 ± 0.7	72.1 ± 1.0	73.9 ± 0.8	$\underline{70.5\pm0.9}$	$\underline{72.9\pm0.8}$	$\underline{76.2 \pm 1.3}$	$\underline{79.4\pm0.9}$	7.1
Full Graph	GCN GAT BGRL GGD	$\begin{array}{c} 54.6 \pm 4.0 \\ 55.1 \pm 4.2 \\ 55.1 \pm 5.8 \\ 58.6 \pm 9.5 \end{array}$	$\begin{array}{c} 54.1 \pm 4.7 \\ 55.3 \pm 5.9 \\ 52.9 \pm 3.1 \\ 56.0 \pm 8.6 \end{array}$	$\begin{array}{c} 50.1 \pm 0.2 \\ 50.1 \pm 0.0 \\ 50.3 \pm 0.2 \\ 50.3 \pm 0.2 \end{array}$	$\begin{array}{c} 50.2 \pm 0.1 \\ 50.0 \pm 0.0 \\ 50.4 \pm 0.2 \\ 50.3 \pm 0.2 \end{array}$	$\begin{array}{c} 53.5 \pm 0.3 \\ \text{OOM} \\ 53.2 \pm 0.3 \\ 53.0 \pm 0.3 \end{array}$	$\begin{array}{c} 53.8 \pm 0.4 \\ \text{OOM} \\ 53.3 \pm 0.3 \\ 53.1 \pm 0.3 \end{array}$	OOM OOM OOM OOM	OOM OOM OOM OOM	OOM OOM OOM OOM	OOM OOM OOM OOM	$11.1 \\ 11.6 \\ 10.9 \\ 10.1$
Full Hypergraph	HGNN UniGCNII ALLSET HCL	$\begin{array}{c} 56.9 \pm 5.3 \\ 55.2 \pm 4.6 \\ 54.2 \pm 4.1 \\ \underline{63.6 \pm 6.9} \end{array}$	$\begin{array}{l} 54.9 \pm 6.3 \\ 52.9 \pm 4.5 \\ 56.1 \pm 5.4 \\ 61.4 \pm 7.9 \end{array}$	$\begin{array}{l} 54.8 \pm 3.7 \\ 51.6 \pm 1.0 \\ 51.0 \pm 0.7 \\ 58.6 \pm 2.6 \end{array}$	$\begin{array}{l} 54.8 \pm 3.6 \\ 51.1 \pm 0.9 \\ 51.7 \pm 0.7 \\ 58.8 \pm 4.2 \end{array}$	$\begin{array}{c} \frac{78.2\pm1.0}{76.7\pm1.2}\\ 60.0\pm2.2\\ \textbf{78.5}\pm\textbf{1.0} \end{array}$	$\frac{80.1 \pm 0.7}{79.0 \pm 1.1} \\ 61.8 \pm 2.5 \\ \mathbf{80.9 \pm 0.7}$	OOM OOM OOM	OOM OOM OOM	OOM OOM OOM	OOM OOM OOM	$6.3 \\ 7.8 \\ 8.6 \\ 4.8$
Partitioned Graph	GCN GAT BGRL GGD	$\begin{array}{c} 51.9 \pm 0.4 \\ 52.4 \pm 0.6 \\ 58.8 \pm 7.7 \\ \textbf{67.7} \pm \textbf{12.7} \end{array}$	$\begin{array}{c} 52.2 \pm 0.4 \\ 53.1 \pm 0.6 \\ 56.8 \pm 5.6 \\ \textbf{67.2} \pm \textbf{13.3} \end{array}$	$50.1 \pm 0.0 \\ 50.4 \pm 0.7 \\ 61.2 \pm 3.9 \\ \underline{59.4 \pm 3.8}$	$50.1 \pm 0.0 \\ 50.5 \pm 1.0 \\ 61.4 \pm 3.9 \\ \underline{59.4 \pm 3.5}$	$\begin{array}{c} 51.6 \pm 0.4 \\ 54.1 \pm 0.7 \\ 65.6 \pm 1.0 \\ 64.8 \pm 1.0 \end{array}$	$\begin{array}{l} 51.7\pm0.4\\ 54.3\pm0.8\\ 66.4\pm0.8\\ 65.5\pm0.6\end{array}$	$\begin{array}{l} 54.9\pm0.7\\ 55.2\pm0.4\\ 65.8\pm0.5\\ 62.8\pm0.8\end{array}$	$\begin{array}{c} 54.3 \pm 0.7 \\ 54.6 \pm 0.4 \\ 67.1 \pm 0.6 \\ 64.3 \pm 0.8 \end{array}$	$\begin{array}{c} {\rm OOM} \\ {\rm OOM} \\ 71.0 \pm 1.0 \\ 68.5 \pm 0.9 \end{array}$	$\begin{array}{c} {\rm OOM} \\ {\rm OOM} \\ {\rm 72.7 \pm 0.9} \\ {\rm 70.5 \pm 0.9} \end{array}$	$ \begin{array}{r} 12.7 \\ 10.3 \\ \underline{3.6} \\ 3.8 \end{array} $
Partitioned Hypergraph	HGNN UniGCNII ALLSET	55.2 ± 6.5 52.6 ± 1.7 53.2 ± 1.6	$\begin{array}{c} 55.2 \pm 7.3 \\ 52.3 \pm 0.9 \\ 54.6 \pm 2.1 \end{array}$	$\begin{array}{c} 52.0\pm1.0\\ 50.4\pm0.3\\ 51.5\pm0.3\end{array}$	$\begin{array}{c} 52.4 \pm 2.0 \\ 50.1 \pm 0.2 \\ 52.5 \pm 0.5 \end{array}$	$\begin{array}{c} 68.8 \pm 1.3 \\ 54.9 \pm 0.6 \\ 57.5 \pm 1.9 \end{array}$	$\begin{array}{c} 69.9 \pm 1.6 \\ 54.6 \pm 0.5 \\ 59.2 \pm 2.8 \end{array}$	$\begin{array}{c} 57.6 \pm 1.6 \\ 59.9 \pm 1.2 \\ 58.1 \pm 1.1 \end{array}$	$\begin{array}{c} 57.5 \pm 2.1 \\ 59.6 \pm 1.3 \\ 60.9 \pm 1.4 \end{array}$	$62.8 \pm 3.0 \\ 65.9 \pm 2.2 \\ OOM$	62.5 ± 3.6 66.0 ± 2.2 OOM	6.6 9.8 8.2
Partitioned Hypergraph	$\begin{array}{c} \operatorname{PCL+PINS} \\ \operatorname{(proposed)} \end{array}$	63.2 ± 10.6	64.0 ± 11.5	58.6 ± 1.0	59.2 ± 1.3	77.6 ± 0.8	79.8 ± 0.6	$\textbf{80.2} \pm \textbf{1.7}$	$\textbf{81.6} \pm \textbf{0.8}$	$\textbf{83.1}\pm\textbf{1.4}$	$\textbf{86.1} \pm \textbf{1.5}$	2.1

*bold: best **<u>underline: second-best</u>

Sunwoo Kim

Datasets, Tasks, and Training Methods for Large-scale Hypergraph Learning

Experimental Results

• RQ2) Effectiveness of PINS on Task 2.

*bold: best

**underline: second-best

Experimental Results

• RQ3) Memory consumption of PINS.

	Method	OGBN-MAG	AMiner	MAG	
Average GPU Memory Usage	PCL w/o PINS	2.344	10.573	23.942	0.1%
(GB)	PCL+PINS	2.347	10.575	23.945	More

Experimental Results

• RQ4) Effectiveness of P-IOS on Task 1 and Task 2.

what was	Task	Data	Metric	PCL w/o P-IOS	PCL+P-IOS
		DBLP	AP	87.1 ± 1.9	$\textbf{90.8} \pm \textbf{1.8}$
F-TO2.			AUROC	$\underline{88.3 \pm 2.0}$	$\textbf{92.9} \pm \textbf{1.2}$
	Task_I	Trivago	AP	$\underline{88.2\pm0.9}$	89.4 ± 1.2
	radr-1		AUROC	$\underline{88.9\pm0.7}$	89.5 ± 0.9
		OGBN-MAG	AP	$\underline{94.1\pm0.4}$	$\textbf{94.8}\pm\textbf{0.4}$
6 8 4			AUROC	$\underline{95.1 \pm 0.2}$	$\textbf{96.2}\pm\textbf{0.2}$
	Task-II	DBLP Trivago	AP	62.4 ± 11.5	64.7 ± 11.5
Original hypergraph			AUROC	$\underline{61.1 \pm 9.9}$	$\textbf{62.3} \pm \textbf{12.3}$
A			AP	$\underline{58.6 \pm 1.0}$	$\textbf{59.2} \pm \textbf{1.0}$
			AUROC	$\underline{58.7 \pm 1.0}$	$\textbf{59.3} \pm \textbf{1.0}$
		OCBN MAC	AP	77.3 ± 0.7	$\textbf{78.8} \pm \textbf{0.9}$
		OGDN-MAG	AUROC	$\underline{79.6\pm0.3}$	80.9 ± 0.7
P-IOS Partition 1 P-IOS Partition 2	*h	old: best		**undarlina:	eacand-bast

**<u>underline: second-best</u>

Datasets, Tasks, and Training Methods for Large-scale Hypergraph Learning

*bold: best

Overview

- Proposed Tasks
- Proposed Datasets
- Proposed Training Method
- O Experimental Results
- O Conclusions

Conclusions

In this work, we present...

Tasks: Two pair-level tasks on hypergraphs,

(Hyperedge disambiguation and Local clustering).

Datasets: Two large-scale benchmark hypergraph datasets,

(AMiner and MAG).

Training methods: Partitioning & Contrastive learning-based scalable

training strategy and its additional tools,

(PCL, PINS, and P-IOS).

Materials

Sunwoo Kim

snmsung Research

Thank You For Listening!

Sunwoo Kim*, Dongjin Lee*, Yul Kim, Jungho Park, Taeho Hwang, Kijung Shin