



Learning Heterogeneous Information Networks for Link Prediction

Speaker: Dr. Hongxu Chen

Host: Dr. Yulei Sui

16/10/2020



My basic information

Hongxu CHEN. 陈红旭



Education Background

- 2009.9 – 2013.6, GUET, B.Sc.
- 2015.2 - 2015.12, The University of Queensland, Master of Computer Science.
- 2016.4 - 2019.10, The University of Queensland, Ph.D. in Computer Science.
 - Supervisor: Dr. Hongzhi Yin & Prof. Xue Li

Education Background

- 2019. 11 – now: Postdoc Research Fellow in Network Science Lab, UTS



My basic information

Professional Services:

- Publicity Co-Chair: [BSEC 2020](#)
- Program Committee Member : [ADMA20](#), [ICONIP 2020](#), [CCL 2020](#), [EMERGING](#) 2020, ADMA19, RSBSD2019, ADMA2018, BDASC'18
- Invited Reviewer (conference): KDD20, IJCAI20, SIGIR20, VLDB20, CIKM19, ICDM19, KDD19, VLDB19, SIGIR19, AAAI19, PAKDD19, ICDM8, CIKM18, ICDM17, CIKM17, WISE17, WISE18.
- Invited Journal Reviewer:
 - Editorial Board of Complexity Journal.
 - IEEE Transactions on Knowledge and Data Engineering (TKDE)*
 - WWW Journal*
 - VLDB Journal*
 - IEEE Transactions on Systems, Man and Cybernetics: Systems*
 - Journal of Complexity*
 - ACM Transactions on Data Science.*
 - Journal of Computer Science & Technology.*
 - IEEE Access*



My basic information

Research profile:

Homepage:

<https://sites.google.com/view/hxchen>

Google Scholar:

<https://scholar.google.com.au/citations?user=W3CtDGQAAAJ&hl=en>

Hongxu Chen

The University of Queensland; University of Technology Sydney
Verified email at uq.edu.au · Homepage

Network Embedding · Social Network Analysis · Data Mining · Recommender Systems

CITED BY

TITLE	CITED BY	YEAR
SPTF: A Scalable Probabilistic Tensor Factorization Model for Semantic-Aware Behavior Prediction	73	2017
H Chen, H Yin, X Sun, H Wang, Y Wang, QH Nguyen 2017 IEEE International Conference on Data Mining(CDM'17)		
PME: Projected Metric Embedding on Heterogeneous Networks for Link Prediction	63	2018
H Chen, H Yin, W Wang, H Wang, QH Nguyen, X Li Proceedings of the 24th ACM SIGKDD International Conference on Knowledge ...		
Air-Attentional intention-aware recommender systems	16	2019
T Chen, H Yin, H Chen, R Yan, QH Nguyen, X Li 2019 IEEE 35th International Conference on Data Engineering (ICDE), 304-315		
People opinion topic model: opinion based user clustering in social networks	16	2017
H Chen, H Yin, X Li, M Wang, W Chen, T Chen Proceedings of the 26th International Conference on World Wide Web Companion ...		
TADA: Trend Alignment with Dual-Attention Multi-Task Recurrent Neural Networks for Sales Prediction	14	2018
T Chen, H Yin, H Chen, L Wu, H Wang, X Zhou, X Li 2018 IEEE International Conference on Data Mining(CDM'18)		
Origin-destination matrix prediction via graph convolution: a new perspective of passenger demand modeling	11	2019
Y Wang, H Yin, H Chen, T Wu, J Xu, K Zheng Proceedings of the 25th ACM SIGKDD International Conference on Knowledge ...		
Exploiting centrality information with graph convolutions for network representation learning	11	2019
H Chen, H Yin, T Chen, QH Nguyen, WC Peng, X Li 2019 IEEE 35th International Conference on Data Engineering (ICDE), 590-601		
Infering Substitutable Products with Deep Network Embedding	5	2019

Cited by

All	Since 2015
Citations	216
h-index	7
i10-index	7

Co-authors

Hongzhi Yin	The University of Queensland
Tong Chen	The University of Queensland
Nguyen Quoc Viet Hung	Griffith University
Haofeng Wang	Alibaba DAMO AI Labs
Xiaofeng Zhou	Professor of Computer Science, ...
Xiaochuai Sun 孙超凯	Associate Professor, Xiamen Uni...
Yang Wang 王杨	Professor, Hebei University of Te...

Dr. Hongxu CHEN 陈红旭

ABOUT:

Dr. Hongxu Chen is now working as a Postdoctoral Research Fellow in Network Science Lab at University of Technology Sydney (UTS). Hongxu completed his Ph.D. at The University of Queensland (UQ), Australia, under the supervision of Dr. Hongzhi Yin and Prof. Xue Li In Data Science Research Group, School of Information Tech. and Electrical Engineering (ITEE).

His research interests include:

- Network Science
 - Complex Network Mining
 - Network Embedding
 - Recommender Systems
 - Social Network Modelling and Analytics
- Office: 11.07.112 UTS, 123 Broadway, Ultimo, Sydney.
E-mail: hongxu.chen[at]uts.edu.au
- Professional Services:
- Editorial Board of Complexity Journal.
 - Publicity Co-Chair: [BSEC 2020](#)
 - Program Committee Member: [ADMA20](#), [ICONIP 2020](#), [CCL 2020](#), [EMERGING 2020](#), ADMA19, RSBSD2019, ADMA2018, BDASC'18
 - Invited Reviewer (conference): KDD20, IJCAI20, SIGIR20, VLDB20, CIKM19, ICDM19, KDD19, VLDB19, SIGIR19, AAAI19, PAKDD19, ICDM8, CIKM18, ICDM17, CIKM17, WISE17, WISE18.
 - Invited Journal Reviewer:





My basic information

Research interests:

- Network/Graph Embedding (Representation Learning)
- Heterogenous Information Networks
- Social Network Analytics
- Recommender Systems

Social Boosted Rec
Bip Exploiti
Convoluti
Hongxu Chen[†]
Hongxu Chen[†] Sch

PME: Project

Multi-level Graph Convolutional Networks for Cross-platform Anchor Link Prediction

Hongxu Chen
University of Technology Sydney
hongxu.chen@uts.edu.au

Tong Chen
The University of Queensland
tong.chen@uq.edu.au

Viet Hung Nguyen
Griffith University
Gold Coast, Australia
quocviethung1@gmail.com

Hongzhi Yin^{*}
The University of Queensland
h.yin@uq.edu.au

Bogdan Gabrys
University of Technology Sydney
Bogdan.Gabrys@uts.edu.au

Weiqing Wang
Monash University
Melbourne, VIC, Australia
teresa.wang@monash.edu

Xue Li[†]
The University of Queensland
Brisbane, QLD, Australia
xueli@itee.uq.edu.au

Xiangguo Sun
Southeast University
sunxiangguo@seu.edu.cn

Katarzyna Musial
University of Technology Sydney
Katarzyna.Musial-Gabrys@uts.edu.au



Outline

- Network Embedding
- Link prediction on Heterogenous Information Networks.
 - Projected Metric Embedding (KDD18)
- Anchor link prediction across different platforms.
 - Multi-level GCN (KDD20)

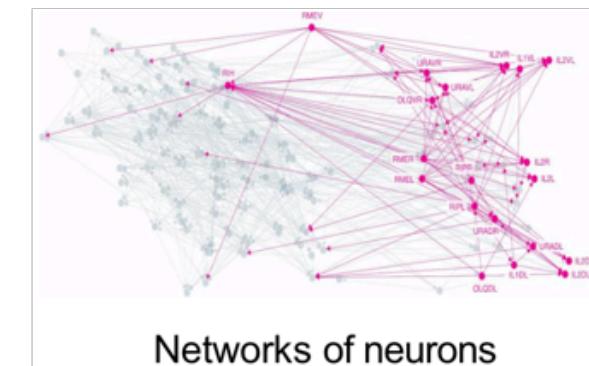
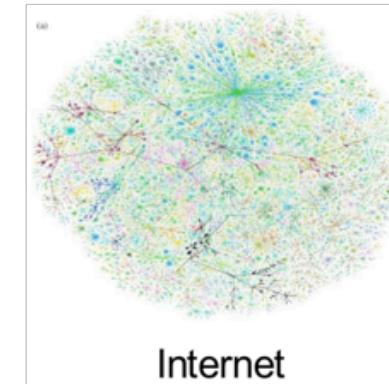
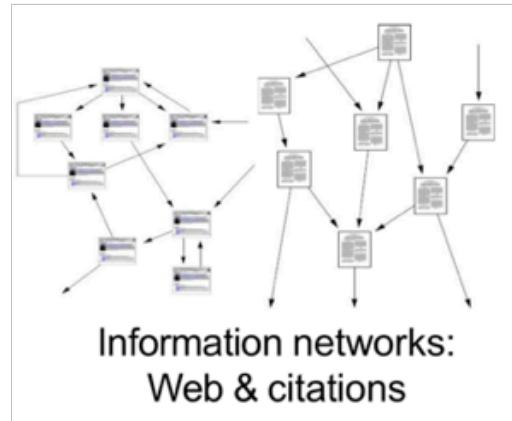
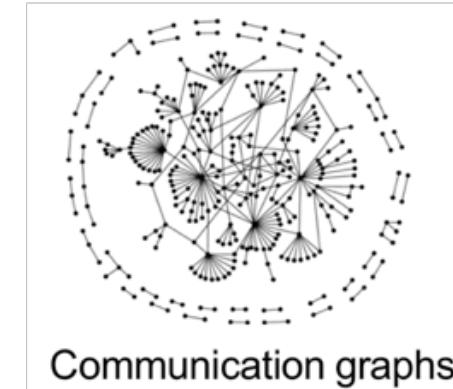


Outline

- Network Embedding
- Link prediction on Heterogenous Information Networks.
 - Projected Metric Embedding (KDD18)
- Anchor link prediction across different platforms.
 - Multi-level GCN (KDD20)



Networks are ubiquitous

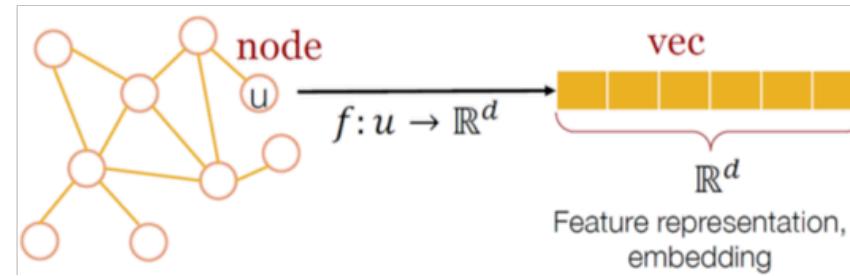




Represent networks by vectors

□ Graph Representation Learning.

- ❖ Also known as Graph Embedding or Network Embedding.
- ❖ Low-dimensional vector for vertices.
- ❖ Effectively preserve network structure.



- ❖ Downstream data mining tasks on graphs:
 - ✓ Link prediction.
 - ✓ Node classification.
 - ✓ Recommendation.
 - ✓ ...



Outline

- Network Embedding
- Link prediction on Heterogenous Information Networks.
 - **Projected Metric Embedding (KDD18)**
- Anchor link prediction across different platforms.
 - Multi-level GCN (KDD20)



Projected Metric Embedding (KDD18)

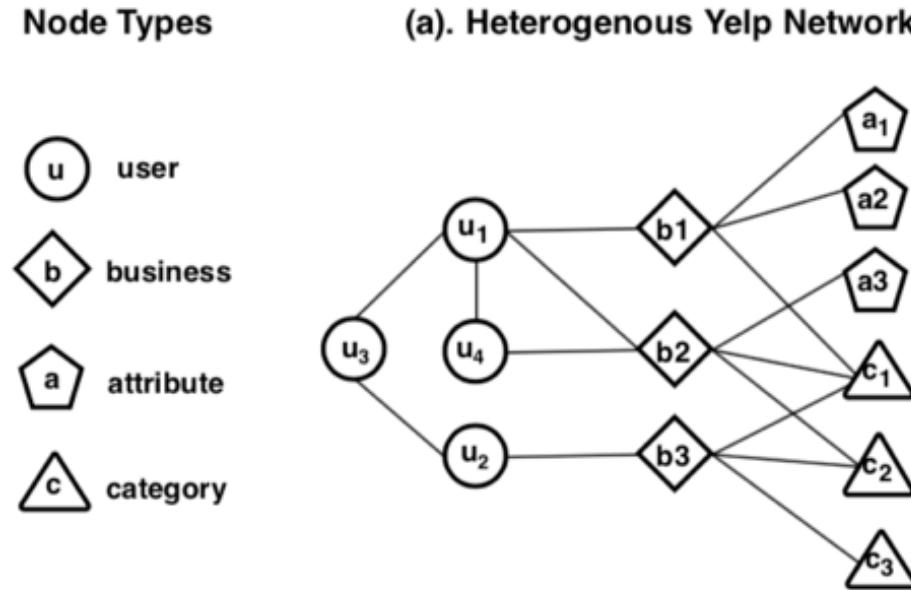


Homogenous Networks

- ❖ Single-Typed Nodes.
- ❖ Single-Typed Links.



Projected Metric Embedding (KDD18)



The Real World: Heterogenous Networks

- ❖ Multiple object types.
- ❖ Multiple link types.

E.g., vertex u_1 is close to both b_2 and u_3 , but these relationships have different semantics. b_2 is a business visited by user u_1 , while u_3 is a friend of u_1 .



Projected Metric Embedding (KDD18)

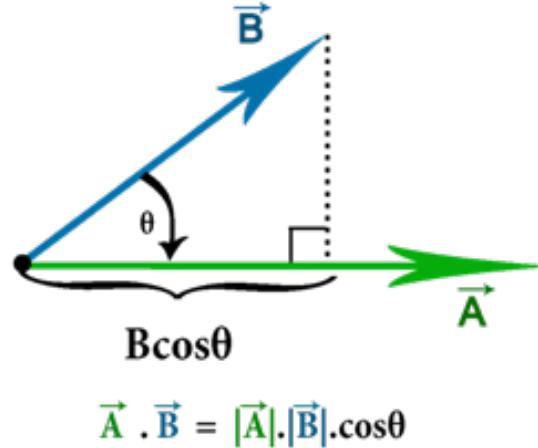
We argue that –

Homogenous networks are information loss projection of heterogenous networks.

- However, most of existing methods only focus on Homogenous Network Embedding that equally treats each type of nodes and links.
- Directly mining information-richer heterogenous networks.
- Compared to homogenous network embedding
The proximity between objects in a HIN is:
 - Not just a measure of closeness or distance.
 - But it is also based on semantics.



Projected Metric Embedding (KDD18)



□ Previous attempts.

To model semantic-specific relationships:

- Metapath2vec (Dong et. al., KDD 2017)
 - EOE (Xu et. al., WSDM 2017)
 - Dot product is used to compute the proximity between different types of nodes
 - Not a metric based distance
 - Violates the crucial triangle inequality
-
- Node A is close to Node C
 - Node B is close to Node C
- } Node A is also close to Node B

Therefore, existing HIN embedding methods (e.g., Metapath2vec and EOE)

- Can only capture local structures (**both A and B are close to C**)
- But fail to capture the second-order proximities. (**A and B are also close**)



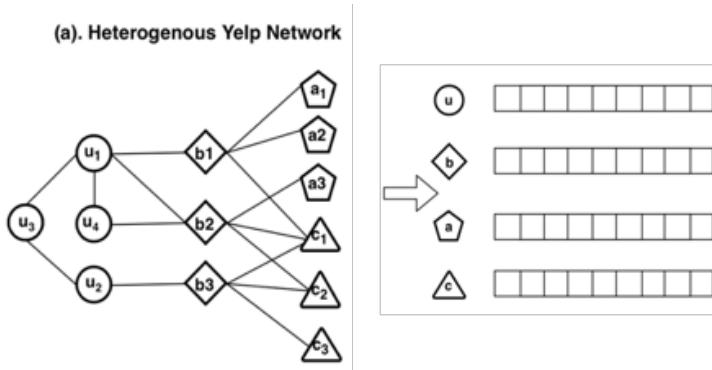
Projected Metric Embedding (KDD18)

Our proposed method

- Projected Metric Embedding for HIN

- ❖ Simultaneously preserves:

- The first-order proximities between nodes
 - And the second-order proximities between nodes



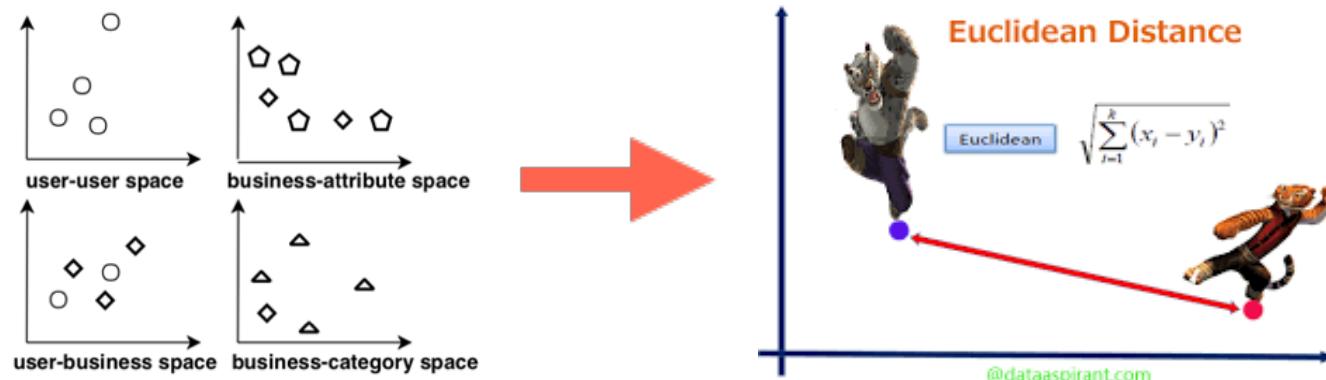
Such that, in the latent space,

$$D(u, v) < D(u, k); e_{uv} \in E, e_{u,k} \notin E$$



Projected Metric Embedding (KDD18)

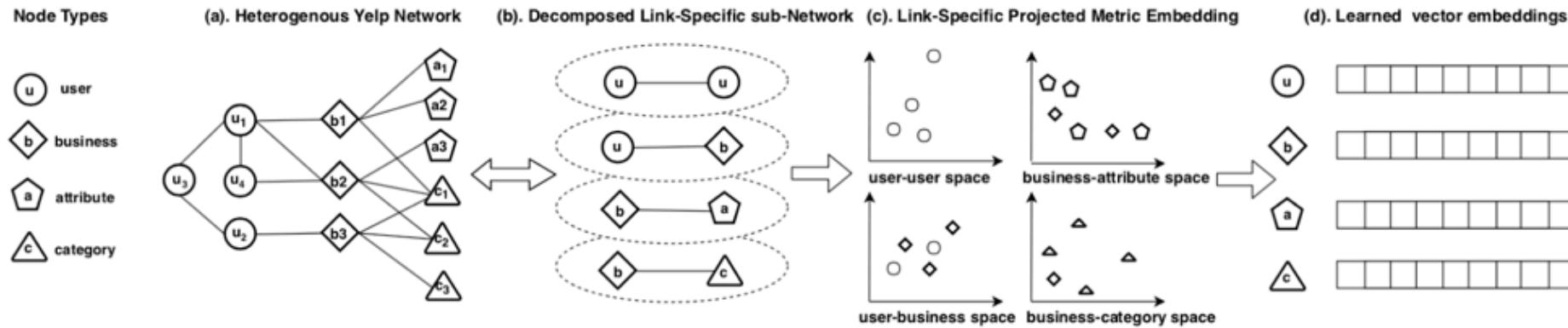
- However, **directly applying the Euclidean distance as a metric will be problematic !**
 - ❖ Mathematically, It is geometrically restrictive and an ill-posed algebraic system.
 - ❖ On the other hand, one object may have multiple aspects.
- To address these issues:
 - ❖ PME introduces relation-specific projection embedding matrices.
 - ❖ Model objects and relations in distinct spaces.
 - One shared object space.
 - Multiple relation spaces. (**relation-specific** object spaces).



Hence, it is possible that some objects are far away from each other in the object space, but are close to each other in the corresponding relation spaces.



Projected Metric Embedding (KDD18)



- ❑ Decompose the HIN to link-specific sub-networks.
- ❑ Metric learning on each link-specific space.

$$v_i^r = M_r v_i \quad d_r(v_i, v_j) = \|M_r v_i - M_r v_j\|, r \in \mathcal{R}$$

- ❑ Learn shared node embeddings and relation-specific space.

$$\begin{aligned} & \min_{v_*, M_*} \sum_{r \in \mathcal{R}} \sum_{(v_i, v_j) \in D_r} \sum_{(v_i, v_k) \notin D_r} [m + f_r(v_i, v_j)^2 - f_r(v_i, v_k)^2]_+ \\ & \text{s.t. } \|v_*\| \leq 1 \quad \text{and} \quad \|M_*\| \leq 1 \end{aligned}$$



Projected Metric Embedding (KDD18)

Model optimization:

- Bi-directional Negative Sampling Strategy

$$\begin{aligned} \min_{\mathbf{v}_*, \mathbf{M}_*} & \sum_{r \in \mathcal{R}} \sum_{(v_i, v_j) \in D_r} \sum_{(v_i, v_k) \notin D_r} [m + f_r(v_i, v_j)^2 - f_r(v_i, v_k)^2]_+ \\ \text{s.t. } & \|\mathbf{v}_*\| \leq 1 \quad \text{and} \quad \|\mathbf{M}_*\| \leq 1 \end{aligned}$$

- Directly optimizing this equation is expensive!
Inspired by the negative sampling techniques

$$\begin{aligned} O = & \sum_{r \in \mathcal{R}} \sum_{(v_i, v_j) \in D_r} \\ & \left(\sum_{k=1}^K E_{v_k} \sim p_n(v) [m + f_r(v_i, v_j)^2 - f_r(v_i, v_k)^2]_+ \right. \\ & \left. + \sum_{k=1}^K E_{v_k} \sim p_n(v) [m + f_r(v_i, v_j)^2 - f_r(v_k, v_j)^2]_+ \right) \end{aligned}$$

first fix vertex v_i and edge type r , then generate K negative vertices v_k

then fix right side of e_{ijr} , and sample K negative vertex from the left side



Projected Metric Embedding (KDD18)

Model optimization:

- Loss-aware Adaptive Positive Sampling Strategy

1. A sequence of the losses for each sub-network.

$$L = (l_1, l_2, l_3, \dots, l_{\|\mathcal{R}\|})$$

2. Simply calculate the sum of the losses.

$$L_{sum} = \sum_{r \in \mathcal{R}} l_r$$

3. Draw a random value within the range of [0,1].

$$x \sim Uniform(0,1)$$

4. To see which interval the random fails into.

$$\left[\sum_{j=0}^{r-1} \frac{l_j}{L_{sum}}, \sum_{j=0}^r \frac{l_j}{L_{sum}} \right)$$

Algorithm 1 Training PME model

Input: A heterogeneous network $G(V, E, W, \mathcal{R})$, number of stochastic gradient steps, N , number of negative samples for each positive sample, K ;
Output: Embeddings for network vertices and relation-specific projection matrix. (i.e., v, M_r);

```
1: iter  $\leftarrow 0$ ;  
2: while iter  $< N$  do  
3:   if iter  $= 0$  then  
4:     Initialize the positive sampling probability as proportional to the original link distribution from  $G$ ;  
5:   else  
6:     Sample  $M$  positive examples based on adaptive positive sampling strategy;  
7:   End if  
8:   For each sampled positive edge, sample  $K$  negative vertices from both sides of the edge;  
9:   Compute gradients and update parameters;  
10:  Censor the norm of  $v$  and projection matrix  $M_r$ ;  
11:  Compute relation-specific subgraph loss, and update the positive sampling probability;  
12:  iter  $\leftarrow$  iter + 1 :  
13: end
```



Projected Metric Embedding (KDD18)

Experimental Results

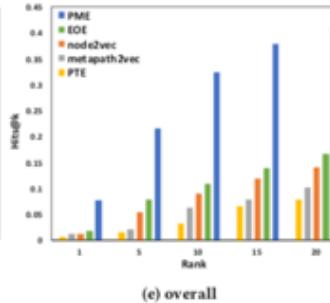
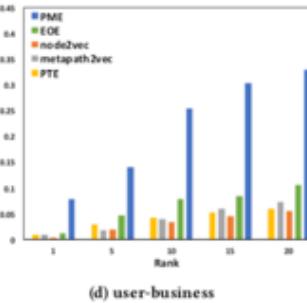
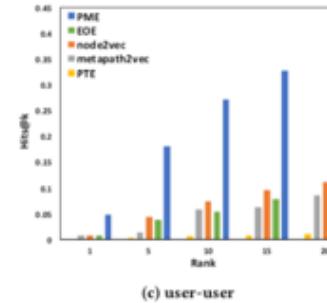
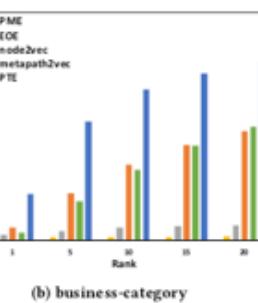
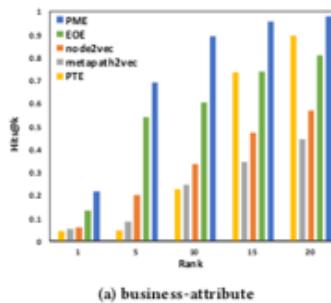
- ❖ Binary Link Classification (Yelp challenge dataset)

Table 4: AUC scores on NV network

	PME	node2vec	PTE	EOE	metapath2vec
Overall	0.9618	0.8789	0.7494	0.8562	0.6232
user-user	0.9672	0.8909	0.6347	0.9033	0.5141
user-business	0.9590	0.8835	0.8615	0.9129	0.8179
business-attribute	0.9376	0.7522	0.8944	0.9201	0.5653
business-category	0.9896	0.9233	0.9652	0.9819	0.7725

- ❖ Prediction accuracy (Hit ratio)

$$Hits@k = \frac{\#hit@k}{\|D_{test}^+\|}$$





Projected Metric Embedding (KDD18)

Code: https://www.dropbox.com/s/o9h7pgkovuryud/KDD18-Hongxu.zip?dl=0&file_subpath=/KDD18-Hongxu

References:

- [1]. Dong, Yuxiao, Nitesh V. Chawla, and Ananthram Swami. "metapath2vec: Scalable representation learning for heterogeneous networks." In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 135-144. 2017.
- [2]. Xu, Linchuan, Xiaokai Wei, Jiannong Cao, and Philip S. Yu. "Embedding of Embedding (EOE) Joint Embedding for Coupled Heterogeneous Networks." In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, pp. 741-749. 2017.



Outline

- Network Embedding
- Link prediction on Heterogenous Information Networks.
 - Projected Metric Embedding (KDD18)
- Anchor link prediction across different platforms.
 - Multi-level GCN (KDD20)



Multi-level Graph Convolutional Networks (KDD20)

Highlights:

- ❑ Cross-platform account matching.
- ❑ Multi-level GCN framework. (Local and Hypergraph levels).
- ❑ Treatments for adapting to large-scale networks.
 - ❖ Network partitioning.
 - ❖ Sub-space reconciliation.

**Multi-level Graph Convolutional Networks for Cross-platform
Anchor Link Prediction**

Hongxu Chen
University of Technology Sydney
hongxuchen@uts.edu.au

Tong Chen
The University of Queensland
tong.chen@uq.edu.au

Hongzhi Yin^{*}
The University of Queensland
h.yin1@uq.edu.au

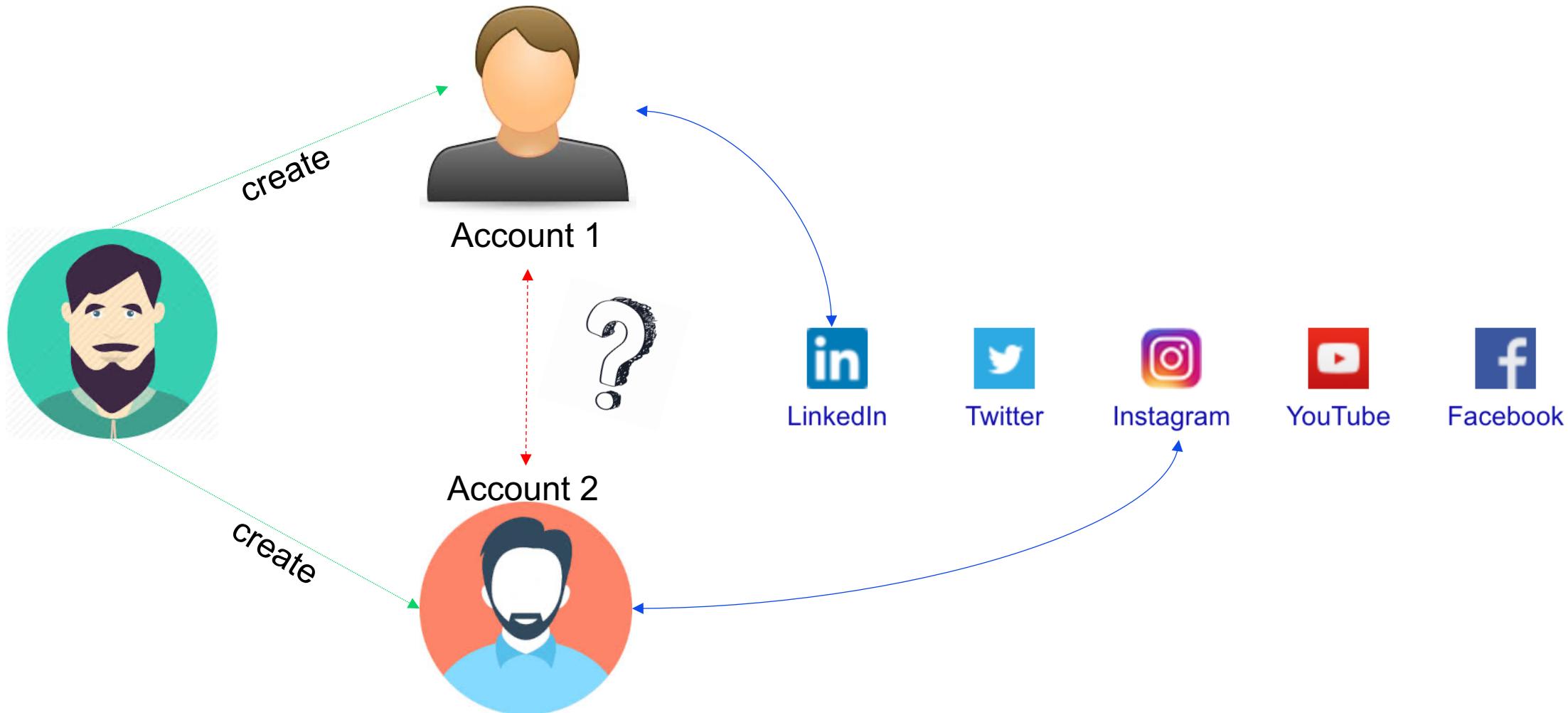
Bogdan Gabrys
University of Technology Sydney
Bogdan.Gabrys@uts.edu.au

Xiangguo Sun
Southeast University
sunxiangguo@seu.edu.cn

Katarzyna Musial
University of Technology Sydney
Katarzyna.Musial-Gabrys@uts.edu.au



Multi-level Graph Convolutional Networks (KDD20)





Multi-level Graph Convolutional Networks (KDD20)

- Cross-platform account matching:
 - ❖ a.k.a. Account Matching, Social Network De-anonymization, Social Anchor Link Prediction.
 - ❖ Beneficial to wide range of applications.

- Current Challenges:
 - ❖ Privacy issue and correctness.
 - Traditional methods by using profiles (self-generated), demographic info (name, pics, location, gender, etc.)
 - ❖ Data-insufficiency.
 - NE based methods need abundant network information.



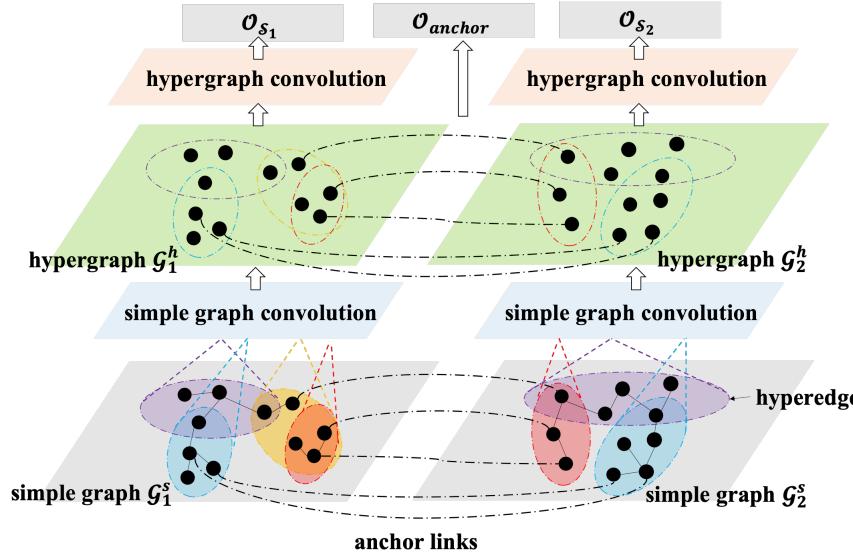
Multi-level Graph Convolutional Networks (KDD20)

- Our method:
 - ❖ Hypergraph integration. (data enrichment)
 - Non-pairwise relations can be captured.
 - ❖ Multi-level GCNs.
 - Simple graph GCN.
 - e.g., friendship, followers.
 - Hypergraph GCN.
 - e.g., N-hop neighbours of a user – **friends circle**.
Centrality-based hypergraphs represent different **social levels**.



Multi-level Graph Convolutional Networks (KDD20)

□ Overview



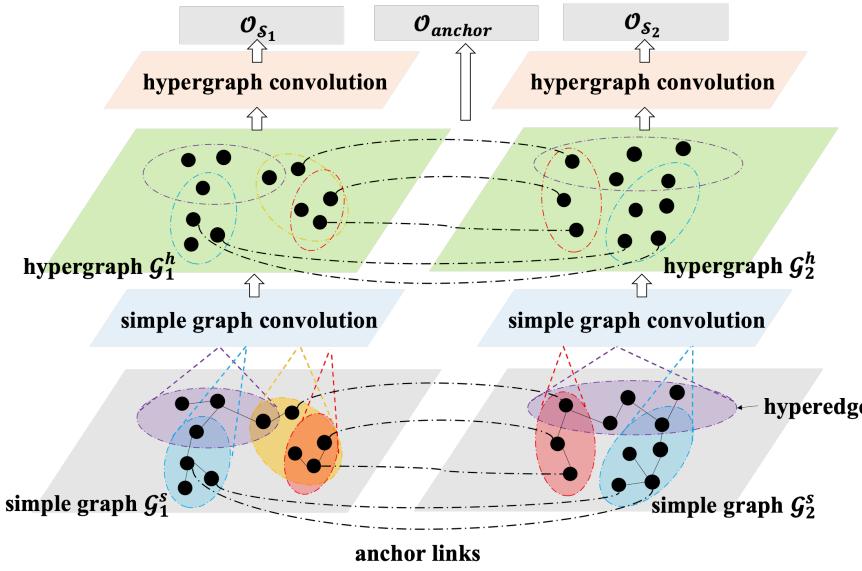
□ Rationale of exploiting and integrating hypergraphs.

- ❖ **hypergraphs provide a more flexible network representation.** - *It can contain additional and richer information compared to individual, single graph GCNs on local network topology.*
- ❖ **GCNs on simple graph are only able to capture the local information.** - *i.e., optimal number of GCN layers is always set to two in most cases because adding more layers cannot significantly improve the performance.*



Multi-level Graph Convolutional Networks (KDD20)

□ Convolution on Simple Graphs:



$$\mathbf{X}_e^{k+1} = \sigma(\mathbf{A}_e \mathbf{X}_e^k \mathbf{W}^k)$$

$$\mathbf{S}_e(v_i, v_j) = \begin{cases} p(v, e), & \text{if } v_i = v_j, v_i \in e \\ 0, & \text{otherwise} \end{cases}$$

$$\mathbf{A}_e = \mathbf{S}_e \hat{\mathbf{A}} \mathbf{S}_e$$

$$\mathbf{X}_{simple}^{k+1} = f(\oplus_{e \in \mathcal{E}^h} \mathbf{X}_e^{k+1})$$

□ Convolution on Hypergraphs:

$$\begin{aligned} \mathbf{X}^{k+1} &= \sigma \left(\left(\mathbf{I}_{|\mathcal{V}|} + \mathbf{D}_n^{-\frac{1}{2}} \mathbf{A}^h \mathbf{D}_n^{-\frac{1}{2}} \right) \mathbf{X}^k \mathbf{W}^k \right) \\ &= \sigma \left(\left(\mathbf{I}_{|\mathcal{V}|} + \mathbf{D}_n^{-\frac{1}{2}} (\mathbf{H} \mathbf{H}^\top - \mathbf{D}_n) \mathbf{D}_n^{-\frac{1}{2}} \right) \mathbf{X}^k \mathbf{W}^k \right) \\ &= \sigma \left(\mathbf{D}_n^{-\frac{1}{2}} \mathbf{H} \mathbf{H}^\top \mathbf{D}_n^{-\frac{1}{2}} \mathbf{X}^k \mathbf{W}^k \right) \end{aligned}$$



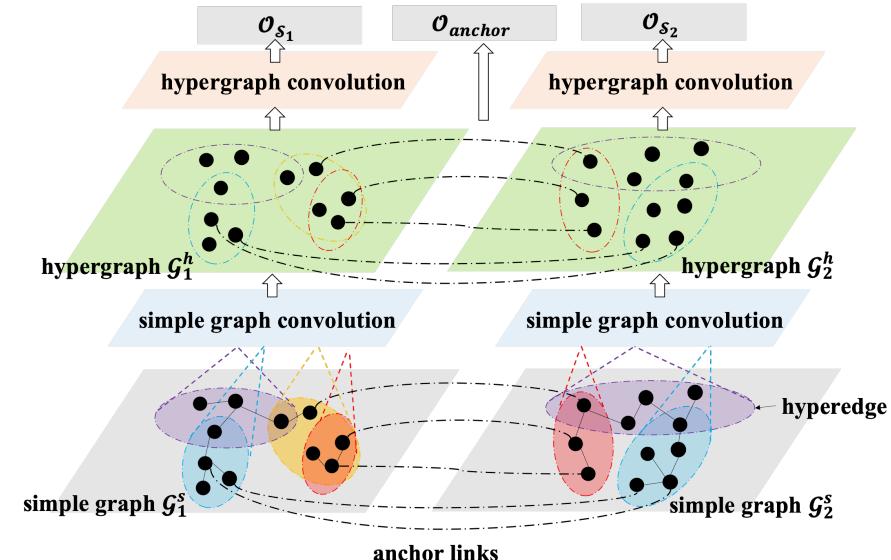
Multi-level Graph Convolutional Networks (KDD20)

□ Learning Network Embedding:

$$\begin{aligned} O_{embedding} &= \sum_{(v_i, v_j) \in \mathcal{E}} \log \eta(\mathbf{x}_i^{K^\top} \mathbf{x}_j^K). \\ &+ \sum_{k=1}^M E_{v_k \propto P(v)} \left[\log (1 - \sigma(\mathbf{x}_i^{K^\top} \mathbf{x}_k^K)) \right] \\ &+ \sum_{k=1}^M E_{v_k \propto P(v)} \left[\log (1 - \sigma(\mathbf{x}_j^{K^\top} \mathbf{x}_k^K)) \right] \end{aligned}$$

□ Anchor Link Prediction:

$$O_{anchor} = \sum_{(v, u) \in \mathcal{S}_{anchor}} \|\mathbf{X}_1^K[v, :] - \phi(\mathbf{X}_2^K[u, :] | \Theta, \mathbf{b})\|^2$$





Multi-level Graph Convolutional Networks (KDD20)

- Handling Large-scale Networks.
 - ❖ Network Partition.
 - ❖ Reconcile Latent Embedding Spaces.

Algorithm 1: Graph Partitioning

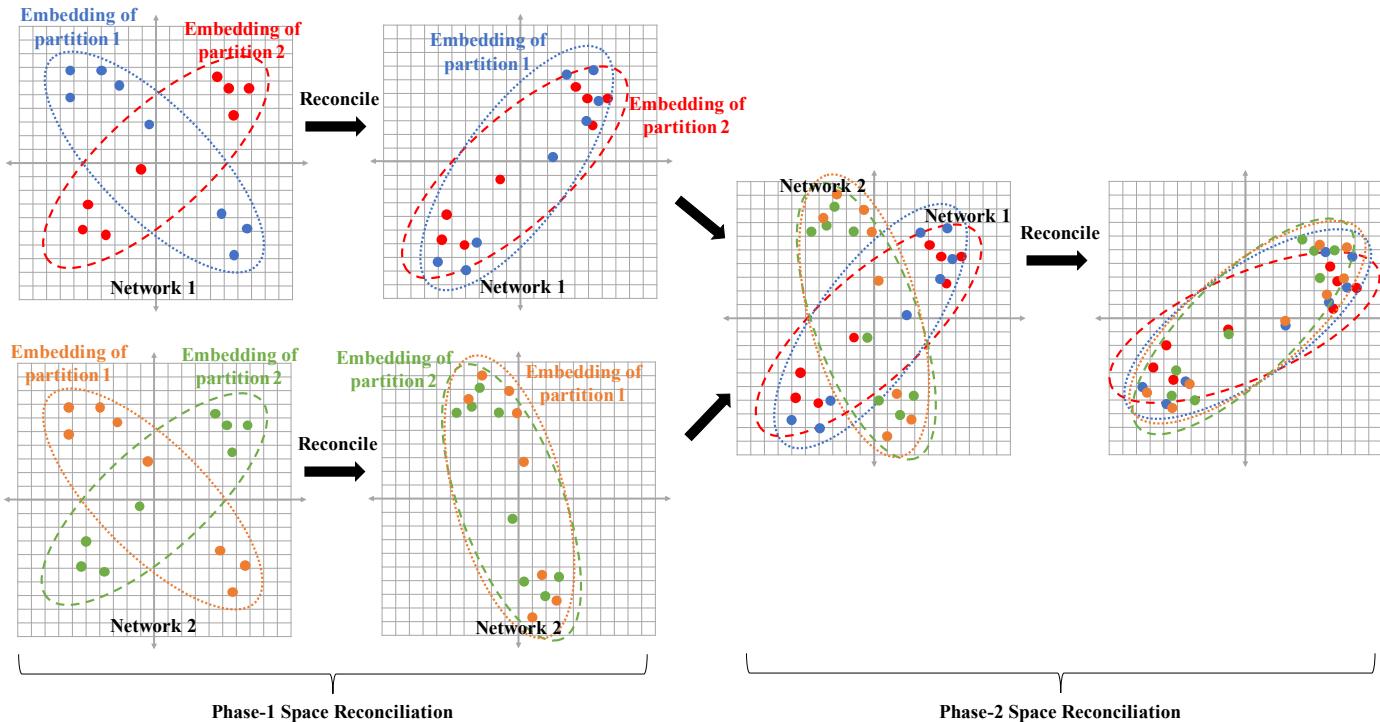
```
Input:  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ ,  $N_{max}$ ,  $N_{min}$ , iteration  $T$ .  
Output: partitions  $P = \{\mathcal{G}_1(\mathcal{V}_1, \mathcal{E}_1), \dots, \mathcal{G}_n(\mathcal{V}_n, \mathcal{E}_n)\}$ .  
1  $P = \text{Louvain}(\mathcal{G})$  //Generating partitions  $P$  from  $\mathcal{G}$   
according to Louvain algorithm[4].  
2 for  $iter$  from 1 to  $T$  do  
3   for partition  $\mathcal{G}' \in P$  do  
4     if  $|\mathcal{V}'| < N_{min}$  then  
5       | add nodes of  $\mathcal{V}'$  into other partitions, delete  $\mathcal{G}'$ .  
6     else if  $N_{min} < |\mathcal{V}'| \leq N_{max}$  then  
7       | continue  
8     else  
9       |  $P_t = \text{Louvain}(\mathcal{G}')$  //Generating partitions  $P_t$   
10      | from  $\mathcal{G}'$  according to Louvain algorithm [4].  
11      |  $P = P \cup P_t$   
12    end  
13  end  
14 return  $P$ 
```

Maximises the modularity using the Louvain algorithm.



Multi-level Graph Convolutional Networks (KDD20)

- Handling Large-scale Networks.
 - ❖ Network Partition.
 - ❖ Reconcile Latent Embedding Spaces.



$$O_{partition} = \sum_{p=2}^P \sum_{v_i \in V_{shared}} \log \sigma \left((f_p(x_i^{(p)}))^T x_i^{(1)} \right)$$



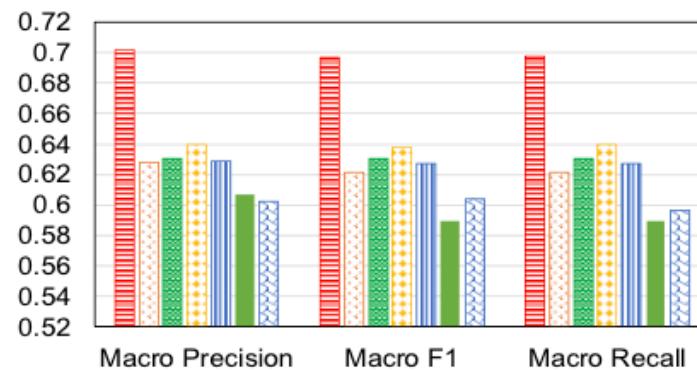
Multi-level Graph Convolutional Networks (KDD20)

- Experiments
 - ❖ Datasets
 - Facebook-Twitter dataset.
 - Douban-Weibo dataset.
 - ❖ Baselines
 - Autoencoder
 - MAH
 - Deepwalk
 - GCN
 - HGCN
 - ❖ Evaluation Metrics:
 - Macro Precession, Recall, and F1. (Treat this task as a binary classification problem)

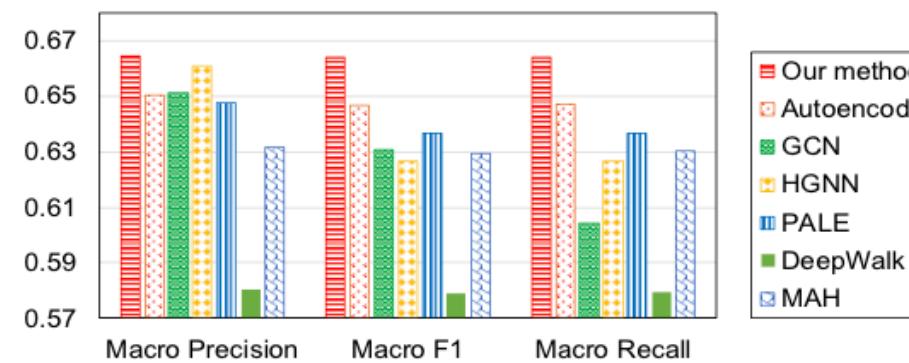


Multi-level Graph Convolutional Networks (KDD20)

- Experiments
 - ❖ Performance on Anchor Link Prediction



(a) Anchor Link Prediction on Facebook-Twitter



(b) Anchor Link Prediction on Douban-Weibo

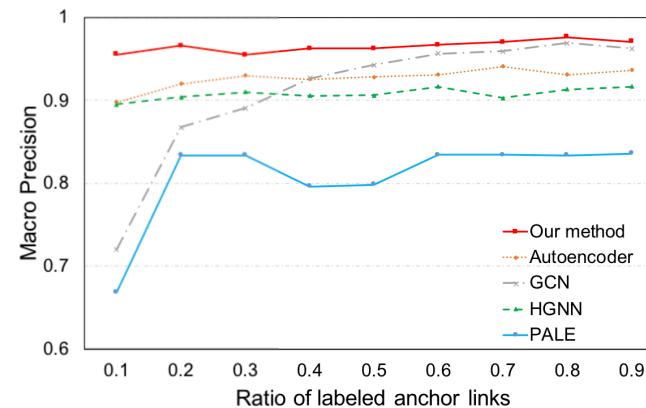


Multi-level Graph Convolutional Networks (KDD20)

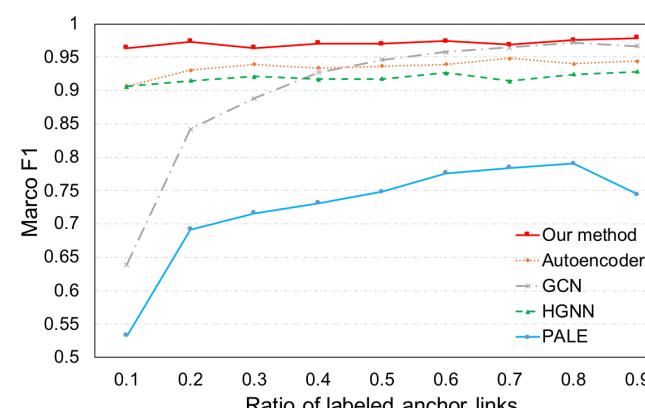
□ Experiments

❖ Analysis on Model Robustness

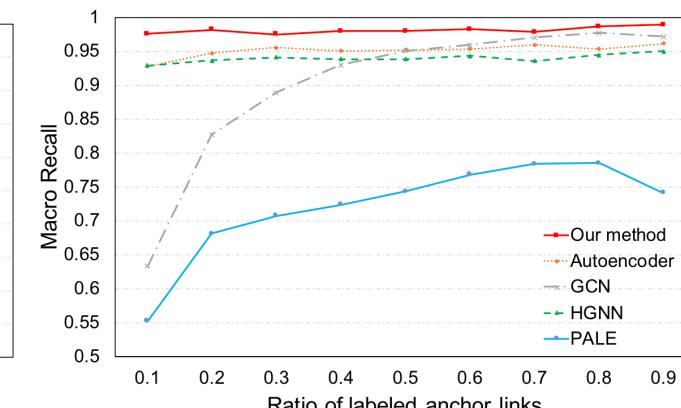
➤ Results w.r.t. observed anchor link percentage.



(a) Macro Precision



(b) Macro F1

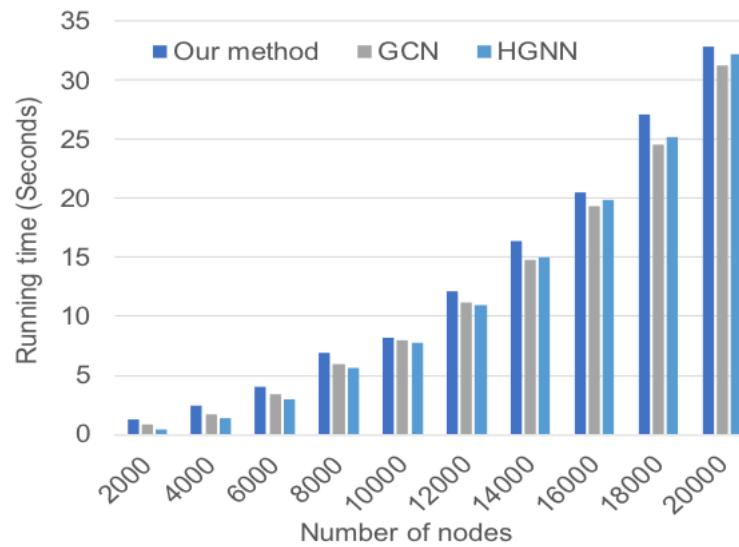


(c) Macro Recall



Multi-level Graph Convolutional Networks (KDD20)

- Experiments
 - ❖ Analysis on Model Efficiency.
 - Forward propagation time w.r.t. network scales.





Multi-level Graph Convolutional Networks (KDD20)

- Preprint paper at:
 - ❖ <https://arxiv.org/pdf/2006.01963.pdf>

- Code available at:
 - ❖ <https://github.com/sunxiangguo/MGCN>



Thanks! Questions?

