





The Power of Summarization in **Network Representation Learning** Danai Koutra Assistant Professor, CSE **Computational Medicine and Bioinformatics (courtesy)**

Joint work with: Caleb Belth, Christos Faloutsos, Brian Gallagher, Aram Galstyan, Mark Heimann, Di Jin, Yike Liu, Ryan Rossi, Tara Safavi, Neil Shah, Chandra Sripada, Pedro Szekely, Yujun Yan, Linhong Zhu, ...

Great Lakes Workshop on Data Science – September 20-22, 2019 In Slides at: https://bit.ly/2m2tlJo







Mark Heimann



Puja Trivedi







Alican Büyükçakır

Yujun Yan







Carol Zheng

1st y.



4th v.

GEMS Lab @ University of Michigan

Research Data and code Lab photos Home People News



Welcome!

We are the Graph Exploration and Mining at Scale (GEMS) lab at the University of Michigan, founded and led by Danai Koutra. Our team researches important data mining and machine learning problems involving interconnected data: in other words, graphs or networks.

From airline flights to traffic routing to neuronal interactions in the brain, graphs are ubiquitous in the real world. Their properties and complexities have long been studied in fields ranging from mathematics to the social sciences. However, many pressing problems involving graph data are still open. One well-known problem is scalability. With continual advances in data generation and storage capabilities, the size of graph datasets has dramatically increased, making scalable graph methods indispensible. Another is the changing nature of data. Real graphs are almost always dynamic, evolving over time. Finally, many important problems in the social and biological sciences involve analyzing not one but *multiple* networks.

So, what do we do?

The problems described above call for principled, practical, and highly scalable graph mining methods, both theoretical and application-oriented. As such, our work c fields like linear algebra, distributed systems, deep learning, and even neuroscie our ongoing projects include:

- Algorithms for multi-network tasks, like matching nodes across networks
- Learning low-dimensional representations of network metric spaces
- Abstracting or "summarizing" a graph with a smaller network
- Analyzing network models of the brain derived from fMRI scans
- Distributed graph methods for iteratively solving linear systems
- Network-theoretical user modeling for various data science applications

We're grateful for funding from Adobe, Amazon, the Army Research Lab, the Mich for Data Science (MIDAS), Microsoft Azure, the National Science Foundation (NSF

Interested?

If you're interested in joining our group, send an email with your interests and CV opportunities@umich.edu.



News

August 2019 2 papers accepted at ICDM

June 2019 1 paper accepted at ASONAM

June 2019 1 paper accepted at PKDD

May 2019 Welcome new PhDs!

May 2019 Tara passes her prelim

April 2019 3 papers accepted to KDD 2019

March 2019 Danai receives NSF CAREER award

January 2019 Danai awarded an Amazon research grant

December 2018













Adobe

Microsoft Azure

A lot of work on network representation learning!

Must-read papers on NRL/NE.

NRL: netv	work representation learning. NE: network embedding.	54. Link Prediction via Subgraph Embedding-Based Cor	nvex Matrix Completion. Zhu Cao, Linlin Wang, Gerard De
Contribut	ted by Cunchao Tu, Yuan Yao and Zhengyan Zhang.	melo. AAAI 2018.	
We releas Represen DeepWall	se OpenNE, an open source toolkit for NE/NRL. This repository <code>r</code> ntation Learning) training and testing framework. Currently, the k, LINE, node2vec, GraRep, TADW and GCN.	55. Generative Adversarial Network based Heterogeneo Citation Recommendation. J. Han, Xiaoyan Cai, Libin	bus Bibliographic Network Representation for Personalized Yang. AAAI 2018.
Survey	papers:	56. DepthLGP: Learning Embeddings of Out-of-Sample Zhu. AAAI 2018. paper	101. Integrative Network Embedding via Deep Joint Reconstruction. Di Jin, Meng Ge, Liang Yang, Dongxiao He, Longbiao Wang, Weixiong Zhang. IJCAI 2018.
1. Repr 2017	resentation Learning on Graphs: Methods and Applications. N . paper	57. Structural Deep Embedding for Hyper-Networks. Ke	102. Scalable Multiplex Network Embedding. Hongming Zhang, Liwei Qiu, Lingling Yi, Yangqiu Song. IJCAI 2018. paper
2. Grap	h Embedding Techniques, Applications, and Performance: A	paper	103. Adversarially Regularized Graph Autoencoder for Graph Embedding. Shirui Pan, Ruiqi Hu, Guodong Long, Jing
3. A Co Zhen	mprehensive Survey of Graph Embedding: Problems, Technic ng, Kevin Chen-Chuan Chang. 2017. paper	58. TIMERS: Error-Bounded SVD Restart on Dynamic Ne Zhu. AAAI 2018. paper	Jiang, Lina Yao, Chengqi Zhang. IJCAI 2018. 104. Dynamic Network Embedding : An Extended Approach for Skip-gram based Network Embedding. Lun Du, Yun
4. Netw	vork Representation Learning: A Survey. Daokun Zhang, Jie Yi	59. Community Detection in Attributed Graphs: An Emb	Wang, Guojie Song, Zhicong Lu, Junshan Wang. IJCAI 2018.
5. A Tu t	torial on Network Embeddings. Haochen Chen, Bryan Perozzi,	Zhang. AAAI 2018.	105. Discrete Network Embedding. Xiaobo Shen, Shirui Pan, Weiwei Liu, Yew-Soon Ong, Quan-Sen Sun. IJCAI 2018.
6. Netw	vork Representation Learning: An Overview.(In Chinese) Cunc	60. Bernoulli Embeddings for Graphs. Vinith Misra, Sumi	106. Deep Attributed Network Embedding. Hongchang Gao, Heng Huang. IJCAI 2018.
2017. 7. Rela t	, paper tional inductive biases, deep learning, and graph networks. P	61. Distance-aware DAG Embedding for Proximity Searc Zhou Zhao, Fanwei Zhu, Kevin Chen-Chuan Chang, M	107. Active Discriminative Network Representation Learning. Li Gao, Hong Yang, Chuan Zhou, Jia Wu, Shirui Pan, Yue Hu. IJCAI 2018.
Santo Kelse Rotvi	n, Alvaro Surfelice, Conzelice, Vinicelas zambalai, Marcese Inalnio oro, Ryan Faulkner, Caglar Gulcehre, Francis Song, Andrew Ballæ ey Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas He linke, Oricl Vinycle, Xuia Li, Bazyan Bassanu, 2018, paper	62. GraphGAN: Graph Representation Learning with Ge Wang, MIAO ZHAO, Weinan Zhang, Fuzheng Zhang, λ	108. ANRL: Attributed Network Representation Learning via Deep Neural Networks. Zhen Zhang, Hongxia Yang, Jiajun Bu, Sheng Zhou, Pinggang Yu, Jianwei Zhang, Martin Ester, Can Wang. IJCAI 2018.
7107-0	unek vinn vinnas rimari kaznan Pastanu zitis taher	63. HARP: Hierarchical Representation Learning for Net AAAI 2018. paper code	109. Feature Hashing for Network Representation Learning. <i>Qixiang Wang, Shanfeng Wang, Maoguo Gong, Yue Wu.</i> IJCAI 2018.
		64. Representation Learning for Scale-free Networks. <i>R</i> 2018. paper	¹ 110. Constructing Narrative Event Evolutionary Graph for Script Event Prediction. <i>Zhongyang Li, Xiao Ding, Ting Liu.</i> IJCAI 2018. paper code
		65. Social Rank Regulated Large-scale Network Embed 2018. paper	111. Deep Inductive Network Representation Learning. Ryan A. Rossi, Rong Zhou, Nesreen K. Ahmed. WWW 2018. paper
			112. A Unified Framework for Community Detection and Network Representation Learning. Cunchao Tu, Xiangkai
	-1		Zeng, Hao Wang, Zhengyan Zhang, Zhiyuan Liu, Maosong Sun, Bo Zhang, Leyu Lin. TKDE 2018. paper



A lot of work on network representation learning!

SDN119 Workshop on We Per Deep Learning DEEP LEARNING DAY



DOOCN-XII: Network Representation Learning

Dynamics On and Of Complex Networks 2019

Frank Room of the UVM Davis Center University of Vermont, Burlington, Vermont, USA Tuesday, May 28th 2019 1:45pm–5:30pm

The Dynamics On and Of Complex Networks (DOOCN) workshop series, aims on exploring statistical dynamics on and of complex networks. *Dynamics on networks* refers to the different types of processes that take place on networks, like spreading, diffusion, and synchronization. Modeling such processes is strongly affected by the topology and temporal variation of the network structure, i.e., by the *dynamics of networks*. Recently, machine learning techniques have been used to model dynamics of massively large complex networks generated from big data, and the various functionalities resulting from the networks. This motivates us to focus on **"Network Representation Learning"** as the significant topic of interest in the 2019 edition.

The First International Workshop on Deep Learning on Graphs: Methods and Applications (DLG'19)

August 5, 2019 Anchorage, Alaska, USA

In Conjunction with the 25th ACM SIGKDD Conference on Knowledge Discovery and Data Mining August 4-8, 2019 Dena'ina Convention Center and William Egan Convention Center Anchorage, Alaska, USA

KDD2019

On Graphs and Manifolds ICLR 2019 Workshop Overview Accepted Papers Schedule Speakers

Program Committee

Organizers

on Learning

ajun

Μι



A lot of work on network representation learning!



inlin Wang, Gerard De

ntation for Personalized

Most work preserves proximity between nodes



aiur

that take place on networks, like spreading, diffusion, and synchronization. Modeling such processes is strongly affected by the topology and temporal variation of the network structure, i.e., by the dynamics of networks. Recently, machine learning techniques have been used to model dynamics of massively large complex networks generated from big data, and the various functionalities resulting from the networks. This motivates us to focus on "Network Representation Learning" as the significant topic of interest in the 2019 edition.

The First International Workshop on Deep Learning on Graphs: Methods and **Applications (DLG'19)**

August 5, 2019 Anchorage, Alaska, USA

In Conjunction with the 25th ACM SIGKDD Conference on Knowledge Discovery and Data Mining August 4-8, 2019 Dena'ina Convention Center and William Egan Convention Center Anchorage, Alaska, USA

on Learning on Graphs and Manifolds

Μι

NRL

Dee

Su



Proximity vs. Structural Similarity





Find similar nodes in the same part of the network

Useful for link prediction, clustering, classification assuming homophily

[Perozzi+ '14; Grover+ '16; Tang+ '15; ...] Find nodes with similar roles all over the network

Useful for role-based classification, transfer learning, ...

[Ribeiro+ '17; Donnat+ '18, ..]

6



What are roles?

- The ways in which nodes / entities / actors relate to each other
- "The behavior expected of a node occupying a specific position" [Homans '67]
 - ♦ e.g., centers of stars
 - members of cliques
 - ♦ peripheral nodes
- Position or equivalence class:

 collection of nodes with the same role



1S LAB [Lorrain & White '71] [Borgatti & Everett '92] [Wasserman & Faust. '94] [Henderson et al. KDD'12] 7

Relevant Sociology Literature

- S.P. Borgatti and M.G. Everett. 1992. Notions of position in social network analysis. Sociological methodology22, 1 (1992)
- Stephen P Borgatti, Martin G Everett, and Jeffrey C Johnson. 2018. Analyzing social networks. Sage
- F. Lorrain and H.C. White. 1971. Structural equivalence of individuals in social networks. Journal of Mathematical Sociology
- S. Boorman, H.C. White: Social Structure from Multiple Networks: II. Role Structures. American Journal of Sociology, 81:1384-1446, 1976.
- R.S. Burt: Positions in Networks. Social Forces, 55:93-122, 1976.
- M.G. Everett, S. P. Borgatti: Regular Equivalence: General Theory. Journal of Mathematical Sociology, 19(1):29-52, 1994.
- K. Faust, A.K. Romney: Does Structure Find Structure? A critique of Burt's Use of Distance as a Measure of Structural Equivalence. Social Networks, 7:77-103, 1985.
- K. Faust, S. Wasserman: Blockmodels: Interpretation and Evaluation. Social Networks, 14:5–61. 1992.
- R.A. Hanneman, M. Riddle: Introduction to Social Network Methods. University of California, Riverside, 2005.
- L.D. Sailer: Structural Equivalence: Meaning and Definition, Computation, and Applications. Social Networks, 1:73-90, 1978.
- M.K. Sparrow: A Linear Algorithm for Computing Automorphic Equivalence Classes: The Numerical Signatures Approach. Social Networks, 15:151-170, 1993.
- S. Wasserman, K. Faust: Social Network Analysis: Methods and Applications. Cambridge University Press, 1994.
- H.C. White, S. A. Boorman, R. L. Breiger: Social Structure from Multiple Networks I. Blockmodels of Roles and Positions. American Journal of Sociology, 81:730-780, 1976.
- D.R. White, K. Reitz: Graph and Semi-Group Homomorphism on Networks and Relations. Social Networks, 5:143-234, 1983.



Sometimes structural similarity is more appropriate than proximity



Graph comparison / classification [KDD'19a; ICDM'19a]





Multiple networks



This talk: Summarization in Network Representation Learning

• Summarization within a GCN for faster training, data denoising and interpretability [ACM KDD'19a]



 Embedding summarization for compression and on-the-fly computation [ACM KDD'19b; PKDD'19]



This talk: Summarization in Network Representation Learning

• Summarization within a GCN for faster training, data denoising and interpretability [ACM KDD'19a]



 Embedding summarization for compression and on-the-fly computation [ACM KDD'19b; PKDD'19]



Interpretable NN-based Classification

- Given a set of networks
 - each associated with a label

- Devise an efficient, interpretable, and parsimonious model
 - that can accurately predict and
 - ♦ explain each label

Interpretable NN-based Classification

Devise an efficient, interpretable, and parsimonious model
 that can accurately predict and
 explain each label phenotype

GEMS LAB



Related Work

- Linear models (PCA, ICA, matrix factorization)
 - + Denoising
 - Fail to capture non-linear interactions
- Neural-network models (different variants of GCN)
 - + Able to model non-linear interactions
 - Need many training samples
 - Need many parameters
 - Long time for training
 - "Black" box

	Fast	Parsimonious	Interpretable
CNN (KDD'17), GraphCNN (NIPS'16)	X	X	X
GCN (ICLR'17), DGCNN (AAAI'18)	\checkmark	X	X
Diffpool (NIPS'18)	\checkmark	Х	inadequate
GroupINN (proposed)	\checkmark	\checkmark	\checkmark



Related Work

- Linear models (PCA, ICA, matrix factorization)
 - + Denoising

Can we build an interpretable NN-based model that is insensitive to noise, parsimonious and able to capture nonlinearities in the prediction task?

pretable

- Long time for training
- "Black" box

GraphCNN (NIPS'16)	\wedge	\wedge	\wedge
GCN (ICLR'17), DGCNN (AAAI'18)	\checkmark	Х	X
Diffpool (NIPS'18)	\checkmark	Х	inadequate
GroupINN (proposed)	\checkmark	\checkmark	\checkmark





CSE











GEMS LAB

Graph Summarization to

- handle noisy data
- train from small samples of high-dim data
- support interpretability







CSE





CSE

1. Node Grouping / Summarization Layer



Recent findings have shown that
 some nodes (ROIs) are most
 related to the phenotype of interest
 → some edges are expected to be more indicative

[Cohen+J Neurosci '16] [Cole+ NeuroImage '07]

• Node grouping layer:

- "hides" the non-indicative edges into a supernode and
- highlights the indicative edges

1. Node Grouping / Summarization Laye

F: learnable common membership matrix •



Real valued importance score of node *i* in the prediction task

Interpretability

- Nonnegative ullet
- Orthogonal (ideally)
- Nodes in supernode **not** • required to be similar / well-connected





2. RWR-based Graph Conv Layer: Intuition



• Random walks:

- useful tool to sample graph structure
- the RWR scores quantify the similarities of other nodes to the seed nodes

• Design: The output \mathbf{Y}_i of layer *i* is: $\mathbf{Y}_i = \sigma(c\mathbf{W}^s\mathbf{Y}_{i-1}\mathbf{Q}_i + \mathbf{I})$





Q1. Comparison with NN-based methods



w/o orthogonality [Ying, NeurIPS' 18] [Wang, KDD' 17] [Kipf, ICLR' 17]





thanks to the summarization layer

GroupINN models are up to 69× faster at training than all the baseline methods, while achieving same or higher accuracy in a variety of prediction tasks.

EMS LAB

Q2. Parsimony of GroupINN

Less is better!

Methods	# parameters	Normalized wrt GroupINN
CNN-1		
CNN-2		
GCN		
Diffpool		
GroupINN	2,892	1×

thanks to the summarization layer

GroupINN can use 15% or much fewer model parameters to achieve comparable or better performance of the baseline methods.



Q3. Interpretability

Acronyms of brain subnetworks. AN: auditory; CBLN: cerebellar; CON: cingulo-opercular; DAN: dorsal attention; FPN: frontoparietal; MRN: memory retrieval; SN: salience; VAN: ventral attention; VN: vision; SM.M: sensory/somatomotor mouth; SM.H: sensory/somatomotor hand

Tasks	Within subnetworks										
	(GroupIN	N		PCA			Diffpool			
Working Memory	MRN	FPN	SN	SM.M	SM.H	AN	SM.M	MRN	CBLN		
Gambling	VAN	VN	DAN	SM.H	AN	SM.M	SM.M	FPN	MRN		
Emotion	SN	CON	VAN	SM.M	SM.H	AN	DMN	MRN	SM.M		
Social	FPN	SN	VAN	SM.M	CBLN	AN	DAN	SM.M	FPN		

- GroupINN finds the most task-positive sub-networks.
- PCA and Diffpool are misled by strong noisy signals from mouth and hand motion.

thanks to the summarization layer

[Cohen, et al., J Neurosci 2016]; [Cole, et al. Neuron 2014];

сем s LAB 🛛 📄 [Yujun Yan, Jiong Zhu, et al. ACM KDD '19] [Davison, et al. PLOS Comp Bio 2015] 27

This talk: Summarization in Network Representation Learning

• Summarization within a GCN for faster training, data denoising and interpretability [ACM KDD'19a]



 Embedding summarization for compression and on-the-fly computation [ACM KDD'19b; PKDD'19]



Embeddings are powerful, but can take up a lot of space!

- For 1B nodes and K=128 \rightarrow 1TB to store the embeddings!
- Can we summarize them?



сем s LAB 📄 [Di Jin, Rossi et al. ACM KDD'19] https://github.com/GemsLab/MultiLENS 29



Graph Summarization Survey

Graph Summarization Methods and Applications: A Survey

YIKE LIU, TARA SAFAVI, ABHILASH DIGHE, and DANAI KOUTRA, University of Michigan Ann Arbor

While advances in computing resources have made processing enormous amounts of data possible, human ability to identify patterns in such data has not scaled accordingly. Efficient computational methods for condensing and simplifying data are thus becoming vital for extracting actionable insights. In particular, while data summarization techniques have been studied extensively, only recently has summarizing interconnected data, or graphs, become popular. This survey is a structured, comprehensive overview of the state-of-the-art methods for summarizing graph data. We first broach the motivation behind and the challenges of graph summarization. We then categorize summarization approaches by the type of graphs taken as ingrup and further organize each category by core methodology. Finally, we discuss applications of summarization on real-world graphs and conclude by describing some open problems in the field.

Additional Key Words and Phrases: Graph mining, g

Yike Liu, Tara Safavi, Abhilash Dighe, and Danai Koutra. 2018. Graph Summarization Methods and Applications: A Survey. ACM Comput. Surv 51, 3, Article 62 (June 2018), 34 pages. https://doi.org/10.1145(31827)



семя LAB 📄 [Liu, Safavi, Dighe, Koutra. ACM Computing Surveys '18.]

Latent Network Summarization



- Given: a graph G(V, E)
- Find: a compressed representation that captures the key structural information and is
 - ♦ independent of graph size (|V|, |E|), and
 - capable of deriving node representations on the fly





Comparison to Related Work

	Input	Repres	ENTATIONS	Meth	HOD	
	Hetero- geneity	Size indep.	Node specific	Proxim. indep.	Scalable	Induc.
Aggregation [2]	✓	×	×	×	1	×
Cosum [34]	X	×	×	\checkmark	×	×
AspEm [31]	\checkmark	×	\checkmark	×	\checkmark	×
metapath2vec [8]	\checkmark	×	\checkmark	×	\checkmark	×
n2vec [11], LINE [32]] X	×	\checkmark	×	\checkmark	×
struc2vec [26]	X	×	\checkmark	\checkmark	×	×
DNGR [6]	X	×	\checkmark	×	×	×
GraphSAGE [12]	\checkmark	×	\checkmark	\checkmark	\checkmark	1
Multi-LENS	\checkmark	\checkmark	\checkmark	\checkmark	~	\checkmark





СSE СSE IAB IAB (Di Jin, Rossi et al. ACM KDD'19) <u>https://github.com/GemsLab/MultiLENS</u> 33

1. Relational functions to aggregate nodewise structural features automatically



сем 📄 [Di Jin, Rossi et al. ACM KDD'19] https://github.com/GemsLab/MultiLENS 34

1. Relational functions to aggregate nodewise structural features automatically



2. Histogram-based heterogeneous

contexts for nodes

[Di Jin, Rossi et al. ACM KDD'19] <u>https://github.com/GemsLab/MultiLENS</u> 35

1. Relational functions to aggregate nodewise structural features automatically



2. Histogram-based heterogeneous contexts for nodes

36

3. Subspace vectors from which we can derive the embeddings

[Di Jin, Rossi et al. ACM KDD'19] https://github.com/GemsLab/MultiLENS



Space comparison

Data	SE	LINE	n2vec	DW	m2vec	AspEm	G2G	ML (MB)
facebook	8.13x	8.48x	12.79x	12.84x	3.82x	8.50x	9.17x	0.58
yahoo	187.1x	180.0x	242.2x	231.0x	79.8x	197.4x	195.8x	0.62
dbpedia	710.0x	714.2x	996.4x	996.2x	-	749.2x	743.6x	0.81
digg	608.2x	612.8x	848.9x	830.3x	259.9x	641.7x	635.2x	0.54
bibson.	1512.1x	1523.0x	2152.5x	2152.5x	-	1595.8x	-	0.75

Multi-LENS requires 4-2152x less output storage space than the other embedding methods.

Data	#Nodes	#Edges	#Node Types	Graph Type	
facebook	4 0 3 9	88 234	1	unweighted	
yahoo-msg	100 058	1057050	2	weighted	
dbpedia	495 936	921 710	4	unweighted	
digg	283 183	4742055	2	unweighted	
bibsonomy	977 914	3 754 828	3	weighted	

[Di Jin, Rossi et al. ACM KDD'19] https://github.com/GemsLab/MultiLENS

Link Prediction



Data	Metric	NA	SE	LINE	DW	n2vec	GR	s2vec	DNGR	m2vec	AspEm	G2G	$\mathbf{ML}(L=1)$	$\mathbf{ML}(L=2)$
facebook	AUC ACC F1 macro	$0.6213 \\ 0.5545 \\ 0.5544$	0.6717 0.5995 0.5716	$0.7948 \\ 0.7210 \\ 0.7210$	$0.7396 \\ 0.6460 \\ 0.6296$	$0.7428 \\ 0.6544 \\ 0.6478$	0.8157 0.7368 0.7367	0.8155 0.7388 0.7387	$0.7894 \\ 0.7062 \\ 0.7060$	$0.7495 \\ 0.7051 \\ 0.7041$	$0.5886 \\ 0.5628 \\ 0.5628$	$0.7968 \\ 0.7274 \\ 0.7273$	0.8703 0.7920* 0.7920 *	0.8709 * 0.7904 0.7905
yahoo-msg	AUC ACC F1 macro	0.7189 0.2811 0.2343	$0.5375 \\ 0.5224 \\ 0.5221$	0.6745 0.6269 0.6265	0.7715 0.6927 0.6897	$0.7830 \\ 0.7036 \\ 0.7016$	$0.7535 \\ 0.6825 \\ 0.6821$	OOT	OOM	$0.6708 \\ 0.6164 \\ 0.6145$	0.5587 0.5379 0.5377	$0.6988 \\ 0.6564 \\ 0.6562$	0.8443 0.7587* 0.7577*	0.8446* 0.7587* 0.7577*
dbpedia	AUC ACC F1 macro	0.6002 0.3998 0.2968	$0.5211 \\ 0.5399 \\ 0.4539$	$0.9632 \\ 0.9111 \\ 0.9110$	$0.8739 \\ 0.8436 \\ 0.8402$	$0.8774 \\ 0.8436 \\ 0.8402$	OOM	OOT	OOM	OOT	$0.6364 \\ 0.5869 \\ 0.5860$	$0.7384 \\ 0.6625 \\ 0.6613$	0.9820 [*] 0.9186 0.9186	$0.9809 \\ 0.9151 \\ 0.9150$
digg	AUC ACC F1 macro	$0.7199 \\ 0.2801 \\ 0.2660$	0.6625 0.6512 0.6223	$0.9405 \\ 0.8709 \\ 0.8709$	0.9664 0.9023 0.9019	$0.9681 \\ 0.9049 \\ 0.9046$	OOM	OOT	OOM	0.9552 0.8891 0.8890	$0.5644 \\ 0.5459 \\ 0.5459$	$0.8978 \\ 0.8492 \\ 0.8492$	0.9894* 0.9596* 0.9595*	0.9893 0.9590 0.9590
bibsonomy	AUC ACC F1 macro	$0.7836 \\ 0.2164 \\ 0.2070$	0.6694 0.6532 0.6064	0.9750 0.9350 0.9349	$0.6172 \\ 0.5814 \\ 0.5781$	0.6173 0.5816 0.5782	OOM	OOT	ООМ	OOT	0.6127 0.5790 0.5772	ООМ	0.9909* 0.9485* 0.9485*	0.9909 0.9466 0.9466

The Multi-LENS node embeddings outperform all the baselines by 3.5–34.3% in AUC.



[Di Jin, Rossi et al. ACM KDD'19] https://github.com/GemsLab/MultiLENS 38

Inductive Anomaly Detection

- Learn summary of G_{t-1} , apply to G_t
- Compute the distance between the embeddings at *t-1* and *t*



🔆 🚰 🛋 📄 [Di Jin, Rossi et al. ACM KDD'19] 🛛 <u>https://github.com/GemsLab/MultiLENS</u>

Can we summarize / compress the embeddings in a different way?





Temporal, Hash-based Node Embeddings

- Given: a time-evolving heterogeneous network G(V, E)
- Learn: a function $\chi: V \rightarrow \{0,1\}^d$ s.t. the derived *d*-dim embeddings
 - 1) preserve similarities in interactions in G,
 - 2) are space-efficient, and
 - accurately capture temporal information and the heterogeneity of the underlying network

[Di Jin, Mark Heimann, et al. PKDD'19] https://github.com/GemsLab/hode2bits

Example: Find similarities in user interactions

• User stitching:

The task of identifying and matching various online references to the same user in real-world web services.

Instance of entity resolution

[Cohen, W.W., Richman, J., KDD 2012] [Dasgupta, A. +, WSDM 2012] [Saha Roy, R.,+ WWW'15] [Kim, Kini, + WWW'17] [Bhattacharya, I., Getoor, L. TKDD 2007], ...





node2bits: Key ideas

- [R1] Graph heterogeneity
 - ♦ General approach that aggregates rich features + node types
- [R2] Temporal dynamics
 - Temporally valid walks to capture short- and long-term interactions
 - Functionally similar nodes are represented by multiple features (structural sim)
- [R3] Efficient similarity comparison
 ♦ Use LSH to hash similar nodes (linear complexity)
- [R4] Low storage requirement

 Binary hashcode with fixed length

[Di Jin, Mark Heimann, et al. PKDD'19] https://github.com/GemsLab/node2bits





node2bits: Workflow



GEMS LAB

[Di Jin, Mark Heimann, et al. PKDD'19] https://github.com/GemsLab/node2bits

Q1: Supervised Identity stitching

short term long term

	Metric	CN	SE	LINE	DW	n2vec	s2vec	DNGR	AspEm	CTDNE	N2B-0	N2B-SH	N2B-LN
bitcoin	AUC ACC F1	$\begin{array}{c} 0.7474 \\ 0.7174 \\ 0.7001 \end{array}$	$\begin{array}{c} 0.5828 \\ 0.5842 \\ 0.5728 \end{array}$	$\begin{array}{c} 0.6071 \\ 0.5842 \\ 0.5828 \end{array}$	$0.6306 \\ 0.6158 \\ 0.6158$	$\begin{array}{c} 0.6462 \\ 0.6158 \\ 0.6157 \end{array}$	$\begin{array}{c} 0.8025 \\ 0.7263 \\ 0.7263 \end{array}$	$0.5909 \\ 0.5526 \\ 0.5525$	$0.5344 \\ 0.5316 \\ 0.5315$	$\begin{array}{c} 0.6987 \\ 0.6000 \\ 0.5964 \end{array}$	$0.7584 \\ 0.7211 \\ 0.7209$	$\begin{array}{c} 0.7609 \\ 0.7268 \\ 0.7271 \end{array}$	$0.7380 \\ 0.6737 \\ 0.6735$
digg	AUC ACC F1	$\begin{array}{c} 0.6217 \\ 0.6217 \\ 0.5585 \end{array}$	$\begin{array}{c} 0.5171 \\ 0.5152 \\ 0.3770 \end{array}$	$0.7878 \\ 0.7694 \\ 0.7683$	$0.7398 \\ 0.6971 \\ 0.6960$	$0.7445 \\ 0.7013 \\ 0.7003$	OOT	OOM	$0.5105 \\ 0.5088 \\ 0.5088$	$0.6967 \\ 0.5915 \\ 0.5884$	0.8185* 0.7982* 0.7958*	$0.7611 \\ 0.7418 \\ 0.7411$	$0.7587 \\ 0.7444 \\ 0.7433$
wiki	AUC ACC F1	0.6997 0.6997 0.6699	OOT	$0.7854 \\ 0.7132 \\ 0.7129$	OOM	OOM	OOT	OOM	$\begin{array}{c} 0.5374 \ 0.5141 \ 0.5141 \end{array}$	$\begin{array}{c} 0.7707 \\ 0.6488 \\ 0.6398 \end{array}$	$0.8230 \\ 0.7145 \\ 0.7088$	0.8259^{*} 0.7510^{*} 0.7476^{*}	$0.8214 \\ 0.7103 \\ 0.7067$
comp-X	AUC ACC F1	$0.5970 \\ 0.5970 \\ 0.5189$	OOM	$0.5000 \\ 0.6757 \\ 0.4032$	OOM	OOM	OOT	OOM	$\begin{array}{c} 0.5213 \ 0.5103 \ 0.5103 \end{array}$	OOM	0.8095^{*} 0.8414^{*} 0.8154^{*}	$0.7496 \\ 0.7959 \\ 0.7581$	$0.7525 \\ 0.7975 \\ 0.7606$

Dynamic + static variants of node2bits outperform baselines by up to 5.2% in AUC and 4.9% in F1 score. Short-term tactic performs better.

[Di Jin, Mark Heimann, et al. PKDD'19] <u>https://github.com/GemsLab/node2bits</u>

Q2: Output storage efficiency



node2bits uses 63-339× less space than the baselines, while achieving comparable or better stitching performance.

GEMS LAB

[Di Jin, Mark Heimann, et al. PKDD'19] https://github.com/GemsLab/node2bits

Summarizing Large Networks: Overview Survey: [CSUR'18]



Beyond Summarization



Structural embeddings for network alignment



[Mark Heimann, Haoming Shen, Tara Safavi, Danai Koutra. ACM CIKM'18]

When is it useful to learn over higherorder networks, and when is it not?



[Caleb Belth, Fahad Kamran, Donna Tjandra, Danai Koutra. IEEE/ACM ASONAM'19]

Take-away messages: Summarization in Network Representation Learning

- Graph summarization Embeddings
- Structural embeddings are less studied, but are more appropriate than proximity-based ones in several tasks
- Summarization within a GCN can help with faster training, data denoising and interpretability [ACM KDD'19a]
- Embedding summarization can achieve compression and on-the-fly computation of representations [ACM KDD'19b; PKDD'19]
- Histograms are powerful at capturing the graph structure [ACM CIKM'18; ACM KDD'19b,c; ECML/PKDD'19; IEEE ICDM'19]
 - ♦ flexible, versatile (heterogeneity, attributes, directionality, weights...),
 - Iess information loss



Talk based on the following papers

- Mark Heimann, Haoming Shen, Tara Safavi, Danai Koutra. REGAL: Representation Learning-based Graph Alignment. ACM CIKM'18.
- Y. Liu, T. Safavi, A. Dighe, D. Koutra. Graph Summarization Methods and Applications: A Survey. ACM Computing Surveys 2018.



- YujunYan, J. Zhu, Marlena Duda, Eric Solarz, Chandra Sripada, Danai Koutra. GroupINN: Groupingbased Interpretable Neural Network-based Classification of Limited, Noisy Brain Data. ACM KDD'19a.
- Di Jin, R. Rossi, Eunyee Koh, Sungchul Kim, Anup. Rao, Danai Koutra. Latent Network Summarization: Bridging Network Embedding and Summarization. ACM KDD'19b.
- D. Jin*, Mark Heimann*, Tara Safavi, Mengdi Wang, Wei Lee, Lindsay Snider, Danai Koutra. Smart Roles: Inferring Professional Roles in Email Networks. ACM KDD'19c.



- Caleb Belth, Fahad Kamran, Donna Tjandra, Danai Koutra. When to Remember Where You Came from: Node Representation Learning in Higher-order Networks. IEEE/ACM ASONAM 2019.
- Tara Safavi, Caleb Belth, Lukas Faber, Davide Mottin, Emmanuel Müller, Danai Koutra. Personalized Knowledge Graph Summarization: From the Cloud to Your Pocket. IEEE ICDM 2019.
- Mark Heimann, Tara Safavi, Danai Koutra. Distribution of Node Embeddings as Multiresolution Features for Graphs. IEEE ICDM 2019.
- Ryan A. Rossi, Di Jin, Sungchul Kim, Nesreen K. Ahmed, Danai Koutra, John Boaz Lee. From Community to Role-based Graph Embeddings. Arxiv 2019.









</>

Mark Heimann Di Jin





Yujun Yan

Thank you! Questions?

http://danaikoutra.com dkoutra@umich.edu

The Power of Summarization in Network Representation Learning

https://github.com/GemsLab/GroupINN

https://github.com/GemsLab/MultiLENS

https://github.com/GemsLab/node2bits

5th y.





















