Advanced Topics in Scalable Learning

CSE 6392 Lecture 2 Graph Neural Network

Junzhou Huang, Ph.D.

Department of Computer Science and Engineering

UT Arlington

Administration

• Course CSE6392

- What: Advanced Topics in Scalable Learning
- When: Friday 1:00 ~ 3:50 pm
- Where: NH 109
- Who: Junzhou Huang (Office ERB 650) jzhuang@uta.edu
- Office Hour: FRIDAY 3:50 ~ 6:00pm and/or appointments
- Webpage: <u>http://ranger.uta.edu/~huang/teaching/CSE6392.htm</u>

(Please check this page regularly)

• Lecturer

- PhD in CS from Rutgers, the State University of New Jersey
- Research areas: machine learning, computer vision, medical image analysis and bioinformatics
- GTA
 - Qifeng Zhou (Office ERB 105B), qxz8706@mavs.uta.edu
 - Office hours: Friday 10:00am ~ 12:00pm and/or appointments
 UT Arlington
 CSE 6392 Advanced Topics in Scalable Learning

Assignment

• Paper Selection

- Each group has two members at most.
- Each group will select at least one paper from the following paper list and then be scheduled to present their selected papers in our class.
- You can choose any papers from the paper lists
- Please talk to the lecturer if you prefer to select a paper out of the list
- The selected paper has to be confirmed by the second week (before the second class)
- GTA will set up the paper selection sheet
- Different groups will present different papers

Grading

• Distribution

- 30% Paper Presentation
- 30% Slide Preparation
- 30% Questions & Answering
- 10% Class Participation
- 100%

• Attention

- No midterm or final exam for this course.
- Please read the selected paper and prepare the final presentation as early as possible
- This is research seminar course. Asking questions and discussion are highly encouraged
- When missing a class due to unavoidable circumstances, PLEASE notify the instructor in advance with any notes/evidences

UT Arlington



Graph is Everywhere

















Graph is Important

• Numerous real-world problems can be summarized as a set of tasks on graphs.



The Power of Deep Learning



CSE 6392 Advanced Topics in Scalable Learning

VGG

GoogleNet ResNet

AlexNet

Deep Learning + Graphs = ?



Deep Graph Learning



The Big Challenges



- Grid-like and sequence-like structure.
- Spital/sequential relations between pixels / units.



Irregular









Large-scale instance

CSE 6392 Advanced Topics in Scalable Learning

VS

The Research Questions



Overview

Foundations

- Preliminary
- The Brief History of Graph Neural Networks

Advances

- Training Deep GNNs
- Scalability of GNNs
- Robustness of GNNs
- Self/Un-Supervised Learning of GNNs
- Other Advanced Topics

Applications

- Social Networks
- Medical Imaging

Preliminaries and Brief History of Graph Neural Networks

UT Arlington

What is the Graph Neural Network?



Graph Neural Network is not a New Thing

Sperduti, Alessandro and Starita, Antonina. 1997

714

IEEE TRANSACTIONS ON NEURAL NETWORKS, VOL. 8, NO. 3, MAY 1997

Supervised Neural Networks for the Classification of Structures

Alessandro Sperduti and Antonina Starita, Member, IEEE

Abstract—Until now neural networks have been used for classifying unstructured patterns and sequences. However, standard neural networks and statistical methods are usually believed to be inadequate when dealing with complex structures because of their feature-based approach. In fact, feature-based approaches usually fail to give satisfactory solutions because of the sensitivity of the approach to the *a priori* selection of the features, and the incapacity to represent any specific information on the relationships among the components of the structures.

	Conce	ptual	Graph
1			

Linear Representation

[HUMAN PROCESS: statement]

-> (AGENT) -> [ACTOR: physician]

Sperduti, Alessandro, and Antonina Starita. "Supervised neural networks for the classification of structures."

A Rapidly Growing Area

ICLR 2020 submissions keyword statistics





https://github.com/shaohua0116/ICLR2020-OpenReviewData

UT Arlington

Preliminaries of Graph Learning

A topological graph



 $\mathcal{G} = {\mathcal{V}, \mathcal{E}}$

which has n nodes and m edges

$$\begin{aligned} \mathcal{V} &= \{v_1, \dots, v_n\} \\ \mathcal{E} &= \{e_1, \dots, e_m\} \end{aligned}$$

Node features: $X \in R^{n \times d}$ Adjacent matrix A: $A_{ij} = 1$, existing edge between i and j $A_{ij} = 0$, not exist edge between i and j

Degree matrix **D** = diag(degree(v_1), ..., degree(v_n))

Laplacian matrix L = D - A

Degree matrix	Adjacency matrix	Laplacian matrix
$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 2 \end{pmatrix}$	$- \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 \end{pmatrix} =$	$\begin{pmatrix} 1 & -1 & 0 & 0 & 0 & 0 \\ -1 & 3 & -1 & 0 & -1 & 0 \\ 0 & -1 & 4 & -1 & -1 & -1 \\ 0 & 0 & -1 & 1 & 0 & 0 \\ 0 & -1 & -1 & 0 & 3 & -1 \\ 0 & 0 & -1 & 0 & -1 & 2 \end{pmatrix}$

The Model of Graph Neural Networks



UT Arlington

The Model of Graph Neural Networks



UT Arlington

GNN 1.0: Understanding GNN as RNN



• The RNN on sequences can be generalized to trees and DAGs.

Sperduti, Alessandro, and Antonina Starita. 1997

UT Arlington

GNN 1.0: Understanding GNN as RNN



Sperduti, Alessandro, and Antonina Starita. 1997 Gori, Marco, Gabriele Monfardini, and Franco Scarselli. 2005 Scarselli, Franco, et al. 2008 Li, Yujia, et.al. 2015, Tai, Kai Sheng et.al, ,2015

UT Arlington

The Brief History of Graph Neural Networks



UT Arlington

GNN 2.0: Understanding GNN as Convolution





- How to perform the convolutions on graphs?
 - Irregular structures.
 - Weighted edges.
 - No orientation or ordering (in general).



GNN 2.0: Understanding GNN as Convolution



UT Arlington

Graph Signal Processing



"Frequency" or "Smoothness" of the signal h





Low frequency graph signal

High frequency graph signal

Eigen decomposition of graph Laplacian $\boldsymbol{L} = \boldsymbol{U} \boldsymbol{\Lambda} \boldsymbol{U}^T$

$$\mathbf{L} = \begin{bmatrix} | & | \\ \mathbf{u}_0 & \cdots & \mathbf{u}_{N-1} \\ | & | \end{bmatrix} \begin{bmatrix} \lambda_0 & 0 \\ & \ddots & \\ 0 & \lambda_{N-1} \end{bmatrix} \begin{bmatrix} - & \mathbf{u}_0 & - \\ & \vdots \\ - & \mathbf{u}_{N-1} & - \end{bmatrix}$$

eigenvalues sorted non-decreasingly: $0 = \lambda_0 \le \lambda_1 \le \dots \le \lambda_{n-1}$

The frequency of an eigenvector of L is its corresponding eigenvalue:

 $u_i^T \boldsymbol{L} u_i = u_i^T \lambda_i u_i = \lambda_i$

UT Arlington

Graph Convolution: Spectral domain \rightarrow Spatial domain

Graph Convolution: input signal x, filter \hat{g} , graph Laplacian L

$$y = x *_{\mathcal{G}} \hat{g} = \boldsymbol{U} \begin{pmatrix} \hat{g}(\lambda_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \hat{g}(\lambda_n) \end{pmatrix} \boldsymbol{U}^T x = \boldsymbol{U} \, \hat{g}(\boldsymbol{\Lambda}) \boldsymbol{U}^T x =: \hat{g}(\boldsymbol{L}) x$$

Parameterization: replace $\hat{g}(\Lambda)$ with $\hat{g}_{\theta} = diag(\theta)$

ChebNet (NeurIPS 2016): parameterize with Chebshev polynomials:

$$y = \hat{g}_{\theta}(L)x = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{L})x \qquad \tilde{L} = \frac{1}{\lambda_{max}} L - I$$

GCN (ICLR 2017): simplified ChebNet K = 1

$$K = 1$$
, suppose $\lambda_{max} = 2$, $\theta \coloneqq \theta_0 = -\theta_1$

 $\hat{g}_{\theta}(\boldsymbol{L})\boldsymbol{x} = \left(\boldsymbol{I} + \boldsymbol{D}^{-\frac{1}{2}}\boldsymbol{A}\boldsymbol{D}^{-\frac{1}{2}}\right)\boldsymbol{x}\boldsymbol{\theta}$

Apply a renormalization trick: $\hat{g}_{\theta}(L)x = \tilde{D}^{-\frac{1}{2}}\tilde{A}\tilde{D}^{-\frac{1}{2}}x\theta$ $\tilde{A} = A + I$ (add self-loop)

UT Arlington

CSE 6392 Advanced Topics in Scalable Learning

Spectral

Spatial

The Brief History of Graph Neural Networks



UT Arlington

Lanczos Network [3]	Graph Wavelet Neural Network [1]	Hyperbolic GCN [2]

[1] Xu, Bingbing, et al. 2018 [2] Chami, Ines, et al. 2019 [3] Liao, Renjie, et al. 2019

UT Arlington



[1] Xu, Bingbing, et al. 2018 [2] Chami, Ines, et al. 2019 [3] Liao, Renjie, et al. 2019

UT Arlington



[1] Xu, Bingbing, et al. 2018 [2] Chami, Ines, et al. 2019 [3] Liao, Renjie, et al. 2019



- Employ Lanczos algorithm to obtain the lowrank approximation of the graph Laplacian $I - \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}.$
- Easy to construct multi-scale Graph Convolution.





$$H_{[:,j]}^{(l+1)} = \sigma\left(\psi_s \sum_{i=1}^p F_{i,j}^{(l)} \psi_s^{-1} H_{[:,i]}^{(l)}\right), \\ j = 1, \dots, q$$

- Use wavelet transform to replace Fourier transform in the original GCN.
- More localized convolution and flexible neighborhood.



Construct the GCN in hyperbolic space.

• Smaller distortion.

2.5

1.5

1.0 0.5

- Suitable for scale-free and hierarchical structure.
- Hyperbolic feature transform. $\boldsymbol{h}_{i}^{(l+1),H} = (\boldsymbol{W}^{(l+1)} \otimes^{K_{l}} \boldsymbol{h}_{i}^{(l),H}) \oplus^{K_{l}} \boldsymbol{b}^{(l+1)}$
- Attention-based hyperbolic aggregation. $y_i^{(l+1),H} = AGG^{K_l}(\mathbf{h}^{(l),H})_i$

[1] Xu, Bingbing, et al. 2018 [2] Chami, Ines, et al. 2019 [3] Liao, Renjie, et al. 2019

UT Arlington

GNN 3.0: GNN with Attention

$$\boldsymbol{h}_{i}^{(l+1)} = \sigma(\sum_{j \in N(v_i)} S_{i,j} \boldsymbol{W}^{(l)} \boldsymbol{h}_{j}^{(l)})$$

$$\boldsymbol{S} = \widetilde{\boldsymbol{D}}^{-\frac{1}{2}} \widetilde{\boldsymbol{A}} \widetilde{\boldsymbol{D}}^{-\frac{1}{2}}$$

Fixed during training



[1] Veličković, Petar, et al. 2018 [2] Zhang, Jiani, et al. 2018 [3] Chang, Heng, et al. 2020

UT Arlington

GNN 3.0: GNN with Attention

The original form:

$$\boldsymbol{h}_{i}^{(l+1)} = \sigma(\sum_{j \in N(v_i)} S_{i,j} \boldsymbol{W}^{(l)} \boldsymbol{h}_{j}^{(l)}) \qquad \boldsymbol{S} = \widetilde{\boldsymbol{D}}$$



[1] Veličković, Petar, et al. 2018 [2] Zhang, Jiani, et al. 2018 [3] Chang, Heng, et al. 2020

UT Arlington

CSE 6392 Advanced Topics in Scalable Learning

Fixed during training

 $\frac{1}{2}\widetilde{A}\widetilde{D}^{-\frac{1}{2}}$

GNN 3.0: GNN with Attention

The original form:

$$\boldsymbol{h}_{i}^{(l+1)} = \sigma(\sum_{j \in N(v_i)} S_{i,j} \boldsymbol{W}^{(l)} \boldsymbol{h}_{j}^{(l)})$$

 $\boldsymbol{S} = \widetilde{\boldsymbol{D}}^{-\frac{1}{2}} \widetilde{\boldsymbol{A}} \widetilde{\boldsymbol{D}}^{-\frac{1}{2}}$

Fixed during training



[1] Veličković, Petar, et al. 2018 [2] Zhang, Jiani, et al. 2018 [3] Chang, Heng, et al. 2020

UT Arlington

Graph Attention Network in Detail



Single head attention



Enrich the model capacity and stabilize the learning process Each head has its own parameters and their outputs can be merged in two ways:

- □ Average

 $\boldsymbol{h}_i^{(l+1)} = \sigma(\sum_{j \in N(v_i)} \boldsymbol{\alpha}_{ij} \boldsymbol{W} \ \boldsymbol{h}_j^{(l)})$

 $\alpha_{ij} = \text{Softmax}(e_{ij})$

$$e_{ij} = \text{LeakyReLU}(\boldsymbol{a}^T (\boldsymbol{W}\boldsymbol{h}_i || \boldsymbol{W}\boldsymbol{h}_j)$$



Attention weights learnt for the Cora dataset
Transformers as GNNs with Multi-head Attention



Chaitanya Joshi. Transformers are graph neural networks, 2020. https://graphdeeplearning.github.io/post/transformers-are-gnns/



- □ Transformers can be viewed as GNNs with multi-head attention as the neighborhood aggregation function
- Transformers for NLP tasks treat the entire sequence as
 the neighborhood

GNN 3.0: GNN with Graph Pooling

Graph Pooling/Coarsening: Convert the node representation to graph representation.

- The most straightforward way: Max/Mean Pooling
- SAGE: Attentive Pooling



Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."

UT Arlington

CSE 6392 Advanced Topics in Scalable Learning

Graph Pooling



Learn the cluster assignment matrix to aggregate the node representations in a hierarchical way.

Incorporate the node features and local structures to obtain a better assignment matrix.

[1] Defferrard, Michaël, et.al. 2016 [2] Ying, Zhitao, et al. 2018 [3] Ma, Yao, et al. 2019

GNN Implementation: Message Passing Framework



Gilmer, Justin, et al. "Neural Message Passing for Quantum Chemistry." ICML. 2017.

UT Arlington

Examples of Message Passing Realizations

Most of current GNNs can be formulated as a message passing process.



Wang, Minjie, et al. "Deep graph library: A graph-centric, highly-performant package for graph neural networks." *arXiv preprint arXiv:1909.01315* (2019). Fey, Matthias and Lenssen, Jan Eric Fast Graph Representation Learning with PyTorch Geometric. (2019). , cite arxiv:1903.02428.

UT Arlington





UT Arlington

Training Deep GNNs

UT Arlington

Training Deep GNNs

Why do we need deep GNNs?
Can GNNs simply go deep?
What impedes GNNs to go deep?
How to make GNNs deep?

The Power of Deep DNNs

- Unprecedented success of deep DNNs in computer vision
- Deep DNNs enable larger receptive fields
- Deep DNNs enable more expressivity





The Power of Deep GNNs

- To GNNs need deep structures to enable larger receptive fields, too?
- What limits the expressive power of GNNs?
 - < The depth *d*
 - The width w
- GNNs significantly lose their power when *capacity*, *dw*, is restricted



Loukas, Andreas. "What graph neural networks cannot learn: depth vs width." International Conference on Learning Representations. 2020.

UT Arlington

The Power of Deep GNNs

The object of the structures to enable larger receptive fields, too?

- What limits the expressive power of GNNs?
 - < The depth *d*
 - The width w

GNNs significantly lose their power when *capacity*, *dw*, is restricted



Loukas, Andreas. "What graph neural networks cannot learn: depth vs width." International Conference on Learning Representations. 2020.

UT Arlington

The Power of Deep GNNs

The boundary of capacity for different problems

problem	bound	problem	bound
cycle detection (odd)	$dw = \Omega(n/{\log n})$	shortest path	$d\sqrt{w} = \Omega(\sqrt{n}/\log n)$
cycle detection (even)	$dw = \Omega(\sqrt{n}/{\log n})$	max. indep. set	$dw = \Omega(n^2/\log^2 n)$ for $w = O(1)$
subgraph verification*	$d\sqrt{w} = \Omega(\sqrt{n}/\log n)$	min. vertex cover	$dw = \Omega(n^2/\log^2 n)$ for $w = O(1)$
min. spanning tree	$d\sqrt{w} = \Omega(\sqrt{n}/\log n)$	perfect coloring	$dw = \Omega(n^2/\log^2 n)$ for $w = O(1)$
min. cut	$d\sqrt{w} = \Omega(\sqrt{n}/\log n)$	girth 2-approx.	$dw = \Omega(\sqrt{n}/\log n)$
diam. computation	$dw = \Omega(n/{\log n})$	diam. ³ /2-approx.	$dw = \Omega(\sqrt{n}/{\log n})$

(Loukas, ICLR'20)

Loukas, Andreas. "What graph neural networks cannot learn: depth vs width." International Conference on Learning Representations. 2020.

UT Arlington

Training Deep GNNs

Why do we need deep GNNs?

- Can GNNs simply go deep?
 - **GCN**: Basic GCN
 - **GraphSAGE**: GCN with improved **aggregation**
 - **ResGCN**: leverage idea from **ResNet**
- **What impedes GNNs to go deep?**
- How to make GNNs deep?

GNNs are Shallow

The set of the set of



Accura	су
--------	----

Citeseer	4 layers	16 layers	64 layers
GCN	76.7	65.2	44.6
GraphSAGE	77.3	72.9	16.9
ResGCN	78.9	78.2	21.2

(Rong et al, ICLR'20)

Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

UT Arlington

Training deep GNNs

- Why do we need deep GNNs?
- Can GNNs simply go deep?
- What impedes GNNs to go deep?
 - Overfitting (Common)
 - Training dynamics (Common)
 - Over-smoothing (Graph Specific)
- How to make GNNs deep?

Training deep GNNs

- Why do we need deep GNNs?
- Can GNNs simply go deep?
- What impedes GNNs to go deep?

Overfitting (Common)

- Training dynamics (Common)
- Over-smoothing (Graph Specific)
- How to make GNNs deep?

Overfitting

GNNs suffer from Overfitting



 $O(dh^2)$

Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

UT Arlington

Training deep GNNs

Why do we need deep GNNs?

Can GNNs simply go deep?

What impedes GNNs to go deep?

• Overfitting (Common)

Training dynamics (Common)

Over-smoothing (Graph Specific)

How to make GNNs deep?

Training dynamics



Training deep GNNs

- Why do we need deep GNNs?
- Can GNNs simply go deep?
- **What impedes GNNs to go deep?**
 - Overfitting (Common)
 - Training dynamics (Common)
 - Over-smoothing (Graph Specific)
- How to make GNNs deep?

- GNNs suffers from over-smoothing
 - Over-smoothing: node representations become less distinguishable with each other when the depth increases

(Li et al. AAAI'18; Chen et al. AAAI'20; Oono et al. ICLR'20; Rong et al. ICLR'20; Huang et al. arXiv'20)



Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

• Over-smoothing also impedes training.



Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

Why GNN works?

- Sy message passing, GNN is able to capture the local structure;
- Several works (Xu et al., 2019, Murphy et al., 2019) show that GNN is equivalent to the Weisfeiler-Lehman (WL) test under a careful design



- When GCNs fail?
 - With linear activation
 - With ReLU activation
 - **With ReLU and bias**



Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

When GCNs fail?

With linear activation

l-step Random Walk

Probability of walking

$$y = \hat{A}^l p_0$$
 where $p_0(i) = 1/d(i)$



Random Walks on Graph

•
$$V_{26} - V_{25} - V_{32} - V_3 - V_{10} \dots$$

•
$$V_5 - V_7 - V_{17} - V_6 - V_{11} \dots$$

• $V_{31} - V_{33} - V_{21} - V_{33} - V_{15}$

Tang, Jian, et al. "Line: Large-scale information network embedding." In WWW, 2015.

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAA/2018.

UT Arlington

When GCNs fail?

With linear activation

l-step Random Walk

$$y = \hat{A}^l p_0$$
 where $p_0(i) = 1/d(i)$

l-layer GCNs

$$Y = \hat{A}^l X W$$

Learnable Probability

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAA/2018.

UT Arlington

When GCNs fail?

With linear activation

l-step Random Walk

$$y = \hat{A}^l p_0$$
 where $p_0(i) = 1/d(i)$

l-layer GCNs

$$Y = \hat{A}^l X W$$

Eigen decomposition

$$Y = \sum_{i=1}^{n} \left(\lambda_{i} u_{i} u_{i}^{\mathsf{T}}\right)^{l} XW$$

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAA/2018.

UT Arlington

Rewrite eigen decomposition

Eigen decomposition

 $(\lambda_1 u_1 u_1^{\mathsf{T}})^l X W + \cdots (\lambda_m u_m u_m^{\mathsf{T}})^l X W + (\lambda_{m+1} u_{m+1} u_{m+1}^{\mathsf{T}})^l X W + \cdots (\lambda_n u_n u_n^{\mathsf{T}})^l X W$

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAA/2018.

UT Arlington

Rewrite eigen decomposition

Eigen decomposition

 $(\lambda_1 u_1 u_1^{\mathsf{T}})^l X W + \cdots (\lambda_m u_m u_m^{\mathsf{T}})^l X W + (\lambda_{m+1} u_{m+1} u_{m+1}^{\mathsf{T}})^l X W + \cdots (\lambda_n u_n u_n^{\mathsf{T}})^l X W$

Suppose graph g has m connected components. It indicates

Eigenvalues

 $1=\lambda_1=\cdots=\lambda_m>\lambda_{m+1}>\cdots>\lambda_n>-1$

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAA/2018.

UT Arlington

Rewrite eigen decomposition



Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAA/2018.

UT Arlington

When
$$l \to +\infty$$
, λ_{m+1} , ..., $\lambda_n \to 0$



The nodes within the same connected component are distinguishable only by their degrees

Li, Qimai, Zhichao Han, and Xiao-Ming Wu. "deep insights into graph convolutional networks for semi-supervised learning." AAA/2018.

UT Arlington

With ReLU activation

Similar to the linear case, but hard to derive the exact

convergence point. Require the notion of subspace (Oono et al., ICLR'20):



Definition 1 (subspace). Let $\mathcal{M} \coloneqq \{ EC | C \in \mathbb{R}^{M \times C} \}$ be an *M*-dimensional subspace in $\mathbb{R}^{N \times C}$, where $E \in \mathbb{R}^{N \times M}$ is orthogonal, i.e. $E^{\mathrm{T}}E = I_M$, and $M \leq N$.

It is proved that (Oono et al., ICLR'20): an infinite-layer GCN will converge to a certain point within a subspace \mathcal{M}



⁽b) Non-Linear Case

Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." ICLR 2020.

UT Arlington



Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." ICLR 2020.

UT Arlington



Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." ICLR 2020.

UT Arlington



Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." ICLR 2020.

UT Arlington

Over-Smoothing of GCNs with bias

• With ReLU and bias

GCNs with bias

$$H_{l+1} = \sigma \left(\hat{A} H_l W_l + \boldsymbol{b}_l \right)$$

Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington
Over-Smoothing of GCNs with bias

 H_L converges to a certain sub-cube $O(\mathcal{M}, r)$ with ReLU and bias



Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

Over-Smoothing of GCNs with bias

With ReLU and bias



Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020





Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

Other methods to measure over-smoothing

One can explicitly measure over-smoothing using distance between node pairs (Chen et al., 2020)



Other metrics see PairNorm (Zhao et al., 2020); GroupNorm (Zhou et al.,

2020)

Chen, et al. "Measuring and Relieving the Over-smoothing Problem for Graph Neural Networks from the Topological View." AAAI 2020 Zhao, Lingxiao, and Leman Akoglu. "PairNorm: Tackling Oversmoothing in GNNs." ICLR. 2020. Zhou et al. "Towards Deeper Graph Neural Networks with Differentiable Group Normalization Kaixiong Zhou Texas A&M University zkxiong@tamu.edu Xiao Huang The Hong." NIPS 2020. UT Arlington

Training deep GNNs

- **Vhy do we need deep GNNs?**
- Can GNNs simply go deep?
- **What impedes GNNs to go deep?**
 - Overfitting (Common)
 - Training dynamics (Common)
 - Over-smoothing (Graph Specific)
- How to make GNNs deep?
 - Architecture refinement
 - Manipulating input (DropEdge)
 - Layer normalizations



There are the term of term

Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

UT Arlington



Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

UT Arlington

Residual connections are helpful (akin to CNNs)

To residual connections alleviate over-smoothing?



Huang, Wenbing, et al. "Tackling Over-Smoothing for General Graph Convolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington



Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington



Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington



Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington



Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

Over-smoothing Layer

For all cases, the over-smoothing speed v^{-1} is controlled by λ_{m+1} (the second-biggest eigenvalue of normalized adjacency matrix)

So how to increase λ_{m+1} ?

Alleviate Over-Smoothing by DropEdge

So how to increase λ_{m+1} ? Drop Edges!

In expectation:



- **The state of the state of the**
- **The gap** $\gamma \mu$ monotonically decreases w.r.t. p;



Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR *2020.* Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

Alleviate Over-Smoothing by DropEdge

 Huang et al., 2020 has considered the re-normalization trick in our analyses, in contrast to Rong et al., 2020

Theorem 1. We denote the original graph as G and the one after dropping certain edges out as G'. Given a small value of ϵ , we assume G and G' will encounter the ϵ -smoothing issue with regard to subspaces \mathcal{M} and \mathcal{M}' , respectively. Then, either of the following inequalities holds after sufficient edges removed.

- The relaxed smoothing layer only increases: $\hat{l}(\mathcal{M}, \epsilon) \leq \hat{l}(\mathcal{M}', \epsilon)$;
- The information loss is decreased: $N dim(\mathcal{M}) > N dim(\mathcal{M}')$.

(Rong et al., 2020)

Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

Alleviate Over-Smoothing by DropEdge

Besides, DropEdge can prevent over-fitting as well!



Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

Alleviate Over-Smoothing by Adjacency Matrix

DropEdge results

Citeseer	4 layers	DropEdges	16 layers	DropEdges	64 layers	DropEdges
GCN	76.7	79.2(+2.5)	65.2	76.8(+11.6)	44.6	45.6(+1.0)
ResGCN	78.9	78.8(-0.1)	78.2	79.4(+1.2)	21.2	75.3(+54.1)
JKNet	79.1	80.2(+1.1)	78.8	80.1(+1.3)	76.7	80.0(+3.3)
IncepGCN	79.5	79.9(+0.4)	78.5	80.2(+1.7)	79.0	79.9(+0.9)
GraphSAGE	77.3	79.2(+1.9)	72.9	74.5(+1.6)	16.9	25.1(+8.2)
APPNP	80.3	80.8(+0.5)	80.2	81.1(+0.9)	80.4	81.3(+0.9)

Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR 2020.

Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

Alleviate Over-Smoothing by Adjacency Matrix

DropEdge results



Rong, Yu, et al. "Dropedge: Towards deep graph convolutional networks on node classification." ICLR *2020.* Huang, Wenbing, et al. "Tackling Over-Smoothing for General GraphConvolutional Networks." arXiv preprint arXiv: 2008.09864, 2020

UT Arlington

Alleviate Over-Smoothing by Weights



Similarly, increasing *s* will also increase v. So how to increase *s*? Increase the initial W_l s.

Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." ICLR 2020.

UT Arlington

Alleviate Over-Smoothing by Weights

Try different **s** as initial



Oono, Kenta, and Taiji Suzuki. "Graph Neural Networks Exponentially Lose Expressive Power for Node Classification." ICLR 2020.

UT Arlington

Training deep GNNs

- **Vhy do we need deep GNNs?**
- Can GNNs simply go deep?
- **What impedes GNNs to go deep?**
 - Overfitting (Common)
 - Training dynamics (Common)
 - Over-smoothing (Graph Specific)

How to make GNNs deep?

- Architecture refinement
- Manipulating input (DropEdge)
- Layer normalizations

Pair Norm: Center and Rescale

PairNorm: Center and rescale (normalize) GCN outputs $\tilde{X} \coloneqq \text{GCN}(A, X)$ to keep the total pairwise squared distance unchanged



Zhao, Lingxiao, and Leman Akoglu. "PairNorm: Tackling Oversmoothing in GNNs." ICLR. 2020.

UT Arlington

Self/Un-Supervised Learning of GNNs

UT Arlington

What we discussed before are supervised



- S. C. P
- Labels are scarce, e.g. molecular property
- Training/Testing tasks are Non I.I.D.

Existing Self-Supervised GNNs

	Node- Classification	Link/Metapath Prediction	Graph- Classification	Graph Reconstruction
Predictive Methods	EP-B [2], GraphSAGE [3], GROVER [7]	S ² GRL[11] SELAR[12]	N-gram Graph [4], PreGNN [5] , GROVER [7] GCC [6]	
Information- based Methods	DGI [8], GMI [9]		InfoGraph [10]	VGAE [1] VRVGA[13] SIG-VAE[14]

[1] Kipf & Welling 2016; [2] Durán & Niepert 2017; [3] Hamilton et al. 2017;
[4] Liu et al. 2019; [5] Hu et al. 2020; [6] Qiu et al. 2020; [7] Rong et al. 2020;
[8] Veli^{*}ckovi[′]c et al. 2019; [9] Peng et al. 2020; [10] Sun et al. 2020
[11] Peng, Zhen, et al. 2020 [12] Hwang, Dasol, et al. 2020
[13] Pan, Shirui, et al. 2018 [14] Hasanzadeh, Arman, et al. 2019

UT Arlington

"In self-supervised learning, the system learns to predict part of its input from other parts of its input." ---- by Yann Lecun

Graphs are highly structured!



Node Classification

• Two typical ways to formulate training loss



I Enforcing Adjacent Similarity

• GraphSAGE (Hamilton et al. 2017)



Enforcing nearby nodes to have similar representations, while enforcing disparate nodes to be distinct:

$$\begin{array}{l} \min - E_{u \sim N(v)} \log(\sigma(h^{T}_{u}h_{v})) - \lambda E_{v_{n} \sim P_{n}(v)} [\log(\sigma(-h^{T}_{v_{n}}h_{v}))] \\ \end{array}$$
Positive Samples Negative Samples

 h_v : representation of target node; h_u : representation of neighbor/positive node; h_{v_n} : representation of disparate/negative node; $P_n(v)$: negative sampling.

UT Arlington

II Reconstruction from neighbors

• EP-B (Durán & Niepert, 2017)





Durán & Niepert. Learning Graph Representations with Embedding Propagation.

UT Arlington

Link/Metapath Prediction

• S²GRL(Peng, Zhen, et al. 2020)



• SELAR(Hwang, Dasol, et al. 2020)



Predict the type of meta path between two nodes

How about graph classification/regression?



N-Gram Graph

• (Liu et al. 2019)



Equivalent to a GNN that needs no train

Liu et al. N-Gram Graph: Simple Unsupervised Representation for Graphs, with Applications to Molecules.

UT Arlington

PreGNN: Node- and Graph-level Pretraining

• (Hu et al. 2020)



Both node- and graph- level training are crucial!

Hu et al. Strategies for Pre-training Graph Neural Networks.

UT Arlington

PreGNN

• (Hu et al. 2020)

Stage I: Node Representation

Enforcing node representation to be similar to its contextual structures:



Hu et al. Strategies for Pre-training Graph Neural Networks.

UT Arlington

PreGNN

• (Hu et al. 2020)

Stage I: Node Representation

Mask random node/edge attribute and predict it, just like Bert:



Context Prediction Attribute Masking

Hu et al. Strategies for Pre-training Graph Neural Networks.

UT Arlington

GCC: Contrastive learning

• (Qiu et al. 2020)



Qiu et al. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training.

UT Arlington
GROVER (Rong et al. 2020)



UT Arlington

• (Rong et al. 2020)



Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.

UT Arlington

• (Rong et al. 2020)



Stage II: Graph-level pretraining

Predicting a graph if contains pre-defined graph motifs.

Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.

UT Arlington

• One more thing: GTransformer



Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.

UT Arlington

We pre-train GROVER with **100 million** parameters on **10 million** unlabeled molecules collected from ZINC15 and Chembl

Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.

UT Arlington

Molecular classification

Classification (Higher is better)								
Dataset	BBBP	SIDER	ClinTox	BACE	Tox21	ToxCast		
# Molecules	2039	1427	1478	1513	7831	8575		
TF_Robust [39]	$0.860_{(0.087)}$	$0.607_{(0.033)}$	$0.765_{(0.085)}$	$0.824_{(0.022)}$	$0.698_{(0.012)}$	$0.585_{(0.031)}$		
GraphConv [23]	$0.877_{(0.036)}$	$0.593_{(0.035)}$	$0.845_{(0.051)}$	$0.854_{(0.011)}$	$0.772_{(0.041)}$	$0.650_{(0.025)}$		
Weave [22]	$0.837_{(0.065)}$	$0.543_{(0.034)}$	$0.823_{(0.023)}$	$0.791_{(0.008)}$	$0.741_{(0.044)}$	$0.678_{(0.024)}$		
SchNet [44]	$0.847_{(0.024)}$	$0.545_{(0.038)}$	$0.717_{(0.042)}$	$0.750_{(0.033)}$	$0.767_{(0.025)}$	$0.679_{(0.021)}$		
MPNN [13]	$0.913_{(0.041)}$	$0.595_{(0.030)}$	$0.879_{(0.054)}$	$0.815_{(0.044)}$	$0.808_{(0.024)}$	$0.691_{(0.013)}$		
DMPNN [61]	$0.919_{(0.030)}$	$0.632_{(0.023)}$	$0.897_{(0.040)}$	$0.852_{(0.053)}$	$0.826_{(0.023)}$	$0.718_{(0.011)}$		
MGCN [29]	$0.850_{(0.064)}$	$0.552_{(0.018)}$	$0.634_{(0.042)}$	$0.734_{(0.030)}$	$0.707_{(0.016)}$	$0.663_{(0.009)}$		
AttentiveFP [59]	$0.908_{(0.050)}$	$0.605_{(0.060)}$	$0.933_{(0.020)}$	$0.863_{(0.015)}$	$0.807_{(0.020)}$	$0.579_{(0.001)}$		
N-GRAM [28]	$0.912_{(0.013)}$	$0.632_{(0.005)}$	$0.855_{(0.037)}$	$0.876_{(0.035)}$	$0.769_{(0.027)}$	_2		
HU. et.al[18]	$0.915_{(0.040)}$	$0.614_{(0.006)}$	$0.762_{(0.058)}$	$0.851_{(0.027)}$	$0.811_{(0.015)}$	$0.714_{(0.019)}$		
GROVER _{base}	$0.936_{(0.008)}$	$0.656_{(0.06)}$	$0.925_{(0.013)}$	$0.878_{(0.016)}$	$0.819_{(0.020)}$	$0.723_{(0.010)}$		
GROVER _{large}	$0.940_{(0.019)}$	$0.658_{(0.023)}$	$0.944_{(0.021)}$	$0.894_{(0.028)}$	$0.831_{(0.025)}$	$0.737_{(0.010)}$		

Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data.

UT Arlington

Existing Self-Supervised GNNs

	Node- Classification	Link/Metapath Prediction	Graph- Classification	Graph Reconstruction
Predictive Methods	EP-B [2], GraphSAGE [3], GROVER [7]	S ² GRL[11] SELAR[12]	N-gram Graph [4], PreGNN [5], GROVER [7] GCC [6]	
Information- based Methods	DGI [8], GMI [9]		InfoGraph [10]	VGAE [1] VRVGA[13] SIG-VAE[14]

[1] Kipf & Welling 2016; [2] Durán & Niepert 2017; [3] Hamilton et al. 2017;

[4] Liu et al. 2019; [5] Hu et al. 2020; [6] Qiu et al. 2020; [7] Rong et al. 2020;

[8] Veličckovičc et al. 2019; [9] Peng et al. 2020; [10] Sun et al. 2020

[11] Hwang, Dasol, et al. 2020 [12] Peng, Zhen, et al.

[13] Pan, Shirui, et al. 2018 [14] Hasanzadeh, Arman, et al. 2019

UT Arlington

What makes a good representation?

Auto-Encoder (AE)



"One natural criterion that we may expect any good representation to meet, at least to some degree, is to retain a significant amount of information about the input." by Vincent et al. 2010

Hinton & Salakhutdinov 2006; Vincent et al. 2010

UT Arlington

Graph Auto-Encoders (VGAE)



Kipf, T. N., & Welling, M. (2016).

UT Arlington

Variational Graph Auto-Encoders (VGAE)



UT Arlington

What makes a good representation?

- A more direct way, other than AE?
 - Yes, Mutual Information (MI).

$$I(X;Y) = D_{KL}(p(X)p(Y)||p(X,Y))$$

= $H(X) - H(X|Y)$
Entropy Conditional Entropy



- $0 \le I(X; Y) \le H(X)$ or H(Y);
- I(X;Y) = 0 iff X and Y are independent random variables;
- I(X;Y) = H(X) = H(Y), if X and Y are determinately related, i.e. H(X|Y) = 0

AE is a lower bound of MI

(Hjelm et al. 2019)



Computing MI is hard and not end-to-end, until recently (CPC, Oord et al., 2018; MINE, Belghazi et al., 2018; Nowozin et al., 2016; Hjelm et al. 2019)

Estimating/Maximizing MI (Hjelm et al. 2019)

Maximize Lower bound of MI

MINE (Belghazi et al., 2018):

$$I^{\text{MINE}}(X;Y) \triangleq E_{p(X,Y)}[T_w(x,y)] - \log E_{p(X)p(Y)}[\exp(T_w(x,y))]$$

JSD MI estimator (Nowozin et al., 2016):

$$I^{\text{JSD}}(X;Y) \triangleq E_{p(X,Y)} \left[\log \sigma \left(T_w(x,y) \right) \right] + E_{p(X)p(Y)} \left[\log \left(1 - \sigma \left(T_w(x,y) \right) \right) \right]$$

InfoNCE MI estimator (Oord et al., 2018):

$$I^{\text{NCE}}(X;Y) \triangleq E_{p(X,Y)}[\log \frac{\exp T_w(x,y)}{\sum_{x' \sim p(X)} \exp T_w(x',y)}]$$

UT Arlington

Deep Graph Infomax (DGI)

• (Velickovic et al. 2019)



The JSD MI estimator is applied:

$$\max_{\text{GNN}} I(X, A; h_i) \approx \max \log(D(h_i; X, A)) + \log(1 - D(\tilde{h}_i; X, A))$$

$$h_i = \text{GNN}(X, A) \qquad \tilde{h}_i \text{ negative sample}$$

UT Arlington

Deep Graph Infomax (DGI)

It is hard to directly compute $D(\tilde{h}_i; X, A)$, thus DGI resorts to readout s = R(X, A):



UT Arlington

Deep Graph Infomax (DGI)

It can be proved that, if the readout s = R(X, A) is injective,

$$\log(D(h_i; \mathbf{s})) + \log(1 - D(\tilde{h}_i; \mathbf{s})) = \log(D(h_i; \mathbf{X}, \mathbf{A})) + \log(1 - D(\tilde{h}_i; \mathbf{X}, \mathbf{A}))$$

It can be also proved that, if |X| = |s| is finite,

$$\max \log(D(h_i; \mathbf{s})) + \log(1 - D(\tilde{h}_i; \mathbf{s})) = \max I(h_i; X, A)$$

UT Arlington

• Some issues in DGI

- Computing MI requires the injectivity of readout function
- ➢ It resorts to graph corruption to generate negative samples
- Distinct encoders and corruption functions for different tasks





• (Peng et al. 2020)



The basic idea of GMI is to compute the MI directly.

Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization.

UT Arlington

We define that,



- \succ It is both feature- and edge- aware;
- ➢ No need to readout or corruption;
- ➢ Feature MI can be further decomposed;



The basic idea of GMI is to compute the MI directly.

Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization.

UT Arlington

• (Peng et al. 2020)

It can be proved that, if certain mild condition meets,

$$I(X;h_i) = \sum_{j \in N(i)} w_{ij} I(x_j; h_i), \text{ for } 0 \leq w_{ij} \leq 1$$



The global MI is decomposed into a weighted sum of local MIs. It is not a bad idea to let $w_{ij} = \sigma(h_i^T h_j)$

We then apply the JSD MI estimator to compute $I(x_j; h_i)$ and $I(\sigma(h_i^T h_j); A_{ij})$

Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization.

UT Arlington

Node Classification

We use a universal backbone (GCN) for all tasks, different from DGI

Algorithm	Transductive			Inductive		
Algorithm	Cora	Citeseer	PubMed	Reddit	PPI	
EP-B loss	79.4 ± 0.1	69.3 ± 0.2	78.6 ± 0.2	93.8 ± 0.03	61.8 ± 0.04	
DGI loss	82.2 ± 0.2	72.2 ± 0.2	78.9 ± 0.3	94.3 ± 0.02	62.3 ± 0.02	
FMI (ours)	78.3 ± 0.1	72.0 ± 0.2	79.1 ± 0.3	94.7 ± 0.03	64.8 ± 0.03	
GMI-mean (ours)	82.7 ± 0.1	73.0 ± 0.3	$\textbf{80.1} \pm \textbf{0.2}$	95.0 ± 0.02	65.0 ± 0.02	
GMI-adaptive (ours)	$\textbf{83.0} \pm \textbf{0.3}$	$\textbf{72.4} \pm \textbf{0.1}$	79.9 ± 0.2	94.9 ± 0.02	64.6 ± 0.03	

Codes: https://github.com/zpeng27/GMI

Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization.

UT Arlington

• Link Prediction

We use an universal backbone (GCN) for all tasks

Algorithm	Cora			BlogCatalog			Flickr			PPI
Aigoritiini	20.0%	50.0%	70.0 %	20.0%	50.0%	70.0%	20.0%	50.0%	70.0%	22.7%
DGI	95.6±0.3	94.6±0.4	94.4±0.2	77.2 ± 0.4	76.4 ± 0.4	75.5 ± 0.3	90.3±0.3	89.0±0.4	74.1±0.7	77.4 ± 0.1
FMI (ours)	97.2±0.2	95.2±0.1	95.0±0.1	81.2 ± 0.2	79.5±0.4	75.1 ± 0.2	92.7±0.3	92.2±0.3	90.6±0.4	79.8±0.2
GMI (ours)	97.9±0.3	96.4±0.2	96.3±0.1	84.1±0.3	83.6±0.2	82.5 ± 0.1	92.0±0.2	90.1±0.3	88.5±0.2	80.0 ± 0.2

Codes: https://github.com/zpeng27/GMI

Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization.

UT Arlington

Summary

	Node- Classification	Link/Metapath Prediction	Graph-Classification	Graph Reconstruction
Predictive Methods	EP-B [2], GraphSAGE [3], GROVER [7]	S ² GRL[11] SELAR[12]	N-gram Graph [4], PreGNN [5] , GROVER [7] GCC [6]	
Information- based Methods	DGI [8], GMI [9]		InfoGraph [10]	VGAE [1] VRVGA[13] SIG-VAE[14]

UT Arlington

Applications

• GNN in Social Networks

GNN in Social Networks

- "Semi-supervised graph classification: A hierarchical graph perspective." WWW 2019
- "Inverse Graph Identification: Can We Identify Node Labels Given Graph Labels?" **arXiv 2020**

Hierarchical Graph Classification

- Hierarchical Graph: A set of graph instances are interconnected via edges.
 - Social network with group structure.
 - Document (graph-of-words) collection with citation relation.



• The Problem: predicts the class label of graph instances in a hierarchical graph.

Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."

UT Arlington

Hierarchical Graph Classification

- The Problem: predicts the class label of graph instances in a hierarchical graph.
- Challenges:
 - How to represent the graphs with arbitrary size into a fixed-length vector?
 - How to incorporate the information of instance level and hierarchical level?



Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."

UT Arlington

Graph Instance Level: Self-Attentive Graph Embedding

- How to represent the graphs with arbitrary size into a fixed-length vector?
- Graph representation learning at different level:
 - Node Level: $G(V, E) \rightarrow H^{n \times v}$
 - Graph Level: $G(V, E) \rightarrow e^{v}$
- SAGE: Self-Attentive Graph Embedding
 - Size invariance ---- Self-attention
 - Permutation invariance ---- GCN Smoothing
 - Node importance ---- Self-attention



Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."

UT Arlington

The Unified Model

- How to incorporate the information of instance level and hierarchical level?
 - Instance Level Model : Graph Level Learning (SEGA)
 - Hierarchical Level Model: Node Level Learning (GCN)
- Feature Sharing: Concatenate the output of SEGA to the input of GCN.
- Disagreement Loss: The disagreement between instance classifier and hierarchical classifier should be minimized.



Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective."

UT Arlington

Applications

• GNN in Medical Imaging

- "Graph CNN for Survival Analysis on Whole Slide Pathological Images", MICCAI 2018
- "Graph Convolutional Nets for Tool Presence Detection in Surgical Videos", IPMI 2019
- "Graph Attention Multi-instance Learning for Accurate Colorectal Cancer Staging", MICCAI 2020

Survival Prediction

- Predict the risk of a certain event occurs.
- Event: part failure, drug adverse reaction or death.
- Application: provides suggestion for clinical interventions

Whole Slide Images

- Large: single WSI size >0.5 GB.
- Complicated: millions of cells.
- Combine local and global features.



UT Arlington

• Cox proportional hazard function

 $\lambda(t|X_i) = \lambda_0(t) \exp(eta_1 X_{i1} + \dots + eta_p X_{ip}) = \lambda_0(t) \exp(X_i \cdot eta)$

• Partial likelihood for event happens on subject *i*:

$$L_i(eta) = rac{\lambda(Y_i|X_i)}{\sum_{j:Y_j \ge Y_i} \lambda(Y_i|X_j)} = rac{\lambda_0(Y_i) heta_i}{\sum_{j:Y_j \ge Y_i} \lambda_0(Y_i) heta_j} = rac{ heta_i}{\sum_{j:Y_j \ge Y_i} heta_j}$$

where, Y is the observation time.

• Join likelihood of all subjects: $L(\beta) =$

$$L(eta) = \prod_{i:C_i=1} L_i(eta)$$

• Log likelihood as object function:

$$\ell(eta) = \sum_{i:C_i=1} \left(X_i \cdot eta - \log \sum_{j:Y_j \geq Y_i} heta_j
ight)$$

UT Arlington

CSE 6392 Advanced Topics in Scalable Learning

 $\theta_i = \exp(X_i \cdot \beta)$



- Pathological Images and Patient Survival Time and Label
 - TCGA, The Cancer Genome Atlas
 - NLST, National Lung Screening Trials

Database	Cancer Subtype	No. Patient	No. WSI	Quality	Avg. Size
TCGA	LUSC	463	535	medium	$0.72~\mathrm{GB}$
TCGA	GBM	365	491	low	$0.50~\mathrm{GB}$
NLST	ADC & SCC	263	425	high	$0.74~\mathrm{GB}$

• Evaluation Metrics- C-index: the fraction of all pairs of patients whose predicted survival times are correctly ordered.

- Yellow regions: high attention values
 - High attention patches : values > 0.9 (attention values (0, 1))





Epoch-5



Epoch-10


GNN in Medical Imaging

Model	LUSC	GBM	NLST
LASSO-Cox [19]	0.5280	0.5574	0.4738
$LASSO-Cox \star$	0.5663	0.5165	0.5663
BoostCI [17]	0.5633	0.5543	0.5705
BoostCI*	0.5800	0.5130	0.5716
EnCox [20]	0.5216	0.5597	0.4883
EnCox*	0.5740	0.5231	0.5742
RSF [12]	0.5066	0.5570	0.5964
$\mathrm{RSF}\star$	0.5492	0.5193	0.5491
MTLSA [16]	0.5386	0.5787	0.6042
$MTLSA \star$	0.5247	0.5630	0.5573
WSISA [21]	0.6380	0.5760	0.6539
GCN-Cox [8]	0.6280	0.5901	0.6845
DeepGraphSurv	0.6606	0.6215	0.7066

* Use our graph features for the survival model.

UT Arlington

CSE 6392 Advanced Topics in Scalable Learning

Future Directions



UT Arlington

CSE 6392 Advanced Topics in Scalable Learning

UT Arlington

CSE 6392 Advanced Topics in Scalable Learning

147

- Sperduti, Alessandro, and Antonina Starita. "Supervised neural networks for the classification of structures." *IEEE Transactions on Neural Networks* 8.3 (1997): 714-735.
- Gori, Marco, Gabriele Monfardini, and Franco Scarselli. "A new model for learning in graph domains." *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*. Vol. 2. IEEE, 2005.
- Scarselli, Franco, et al. "The graph neural network model." IEEE Transactions on Neural Networks 20.1 (2008): 61-80.
- Li, Yujia, et al. "Gated graph sequence neural networks." arXiv preprint arXiv:1511.05493 (2015).
- Defferrard, Michaël, Xavier Bresson, and Pierre Vandergheynst. "Convolutional neural networks on graphs with fast localized spectral filtering." Advances in neural information processing systems. 2016.
- Tai, Kai Sheng, Richard Socher, and Christopher D. Manning. "Improved semantic representations from tree-structured long short-term memory networks." arXiv preprint arXiv:1503.00075 (2015).
- Bruna, Joan, et al. "Spectral networks and locally connected networks on graphs." arXiv preprint arXiv:1312.6203 (2013).
- Niepert, Mathias, Mohamed Ahmed, and Konstantin Kutzkov. "Learning convolutional neural networks for graphs." International conference on machine learning. 2016.
- Kipf, Thomas N., and Max Welling. "Semi-supervised classification with graph convolutional networks." arXiv preprint arXiv:1609.02907 (2016).
- Xu, Bingbing, et al. "Graph Wavelet Neural Network." International Conference on Learning Representations. 2018.
- Chami, Ines, et al. "Hyperbolic graph convolutional neural networks." Advances in neural information processing systems. 2019.
- Liao, Renjie, et al. "LanczosNet: Multi-Scale Deep Graph Convolutional Networks." International Conference on Learning Representations. 2018.
- Veličković, Petar, et al. "Graph Attention Networks." International Conference on Learning Representations. 2018.
- Zhang, Jiani, et al. "Gaan: Gated attention networks for learning on large and spatiotemporal graphs." arXiv preprint arXiv:1803.07294 (2018).
- Chang, Heng, et al. "Spectral Graph Attention Network." arXiv preprint arXiv:2003.07450 (2020).
- Ying, Zhitao, et al. "Hierarchical graph representation learning with differentiable pooling." Advances in neural information processing systems. 2018.
- Ma, Yao, et al. "Graph convolutional networks with eigenpooling." Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2019.
- Atwood, James, and Don Towsley. "Diffusion-convolutional neural networks." Advances in neural information processing systems. 2016.
- Abu-El-Haija, Sami, et al. "MixHop: Higher-Order Graph Convolutional Architectures via Sparsified Neighborhood Mixing." International Conference on Machine Learning. 2019.
- Klicpera, Johannes, Aleksandar Bojchevski, and Stephan Günnemann. "Predict then Propagate: Graph Neural Networks meet Personalized PageRank." International Conference on Learning Representations. 2018.
- Gilmer, Justin, et al. "Neural Message Passing for Quantum Chemistry." ICML. 2017.
- Li, Jia, et al. "Semi-supervised graph classification: A hierarchical graph perspective." The World Wide Web Conference. 2019.

- Hamilton, Will, Zhitao Ying, and Jure Leskovec. "Inductive representation learning on large graphs." Advances in neural information processing systems. 2017.
- Chen, Jianfei, Jun Zhu, and Le Song. "Stochastic Training of Graph Convolutional Networks with Variance Reduction." International Conference on Machine Learning. 2018.
- Chen, Jie, Tengfei Ma, and Cao Xiao. "FastGCN: Fast Learning with Graph Convolutional Networks via Importance Sampling." International Conference on Learning Representations. 2018.
- Huang, Wenbing, et al. "Adaptive sampling towards fast graph representation learning." Advances in neural information processing systems. 2018.
- Chiang, Wei-Lin, et al. "Cluster-GCN: An efficient algorithm for training deep and large graph convolutional networks." Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2019.
- Zeng, Hanqing, et al. "GraphSAINT: Graph Sampling Based Inductive Learning Method." International Conference on Learning Representations. 2020.
- Morris et al. Weisfeiler and Leman Go Neural Higher-order Graph Neural Networks. AAAI, 2019.
- Xu et al. How Powerful Are Graph Neural Networks. ICLR, 2019.
- Maron et al. Invariant and Equivariant Graph Networks. ICLR, 2019.
- Maron et al. On the Universality of Invariant Networks. ICML, 2019.
- Maron et al. Provably Powerful Graph Networks, NeurIPS. 2019.
- Dehmamy et al. Understanding the Representation Power of Graph Neural Networks in Learning Graph Topology. NeurIPS, 2019.
- Sato et al. Approximation Ratios of Graph Neural Networks for Combinatorial Problems. NeurIPS, 2019.
- Loukas, What graph neural networks cannot learn: depth vs width. ICLR, 2020.
- Garg et al. Generalization and Representational Limits of Graph Neural Networks. ICML, 2020.
- Shervashidze et al. Weisfeiler-Lehman Graph Kernels. JRML, 2011.
- Cybenko, G. Approximation by superpositions of a sigmoidal function. Mathematics of control, signals and systems, 2(4):303–314, 1989.
- Hornik, K. Approximation capabilities of multilayer feedforward networks. Neural networks, 4(2):251–257, 1991.
- Chen et al. On the equivalence between graph isomorphism testing and function approximation with GNNs. NeurIPS, 2019.
- Kipf & Welling. Variational Graph Auto-Encoders. arXiv, 2016.
- Durán & Niepert. Learning Graph Representations with Embedding Propagation. NeurIPS, 2017.
- Liu et al. N-Gram Graph: Simple Unsupervised Representation for Graphs, with Applications to Molecules. NeurIPS, 2019.
- Hu et al. Strategies for Pre-training Graph Neural Networks. ICLR, 2020.
- Qiu et al. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training. KDD, 2020.
- Rong et al. GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data. arXiv, 2020.
- Veličković et al. Deep Graph Infomax. ICLR, 2019.
- Peng et al. Graph Representation Learning via Graphical Mutual Information Maximization. WWW, 2020.
- Sun et al. InfoGraph: Unsupervised and Semi-supervised Graph-Level Representation Learning via Mutual Information Maximization. ICLR, 2020.

UT Arlington

CSE 6392 Advanced Topics in Scalable Learning

- Hinton, G. E., & Salakhutdinov, R. R. Reducing the Dimensionality of Data with Neural Networks. Science, 2006.
- Vincent et al. Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion. JMLR, 2010.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. NeurIPS, 2018.
- Belghazi et al. Mine: mutual information neural estimation. ICML, 2018.
- Nowozin et al. f-gan: Training generative neural samplers using variational divergence minimization. NeurIPS, 2016.
- Hjelm et al. Learning deep representations by mutual information estimation and maximization. ICLR, 2019.
- Wang, X., & Gupta, A. Videos as Space-Time Region Graphs. ECCV, 2018.
- Yan et al. Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition. AAAI, 2018.
- Zeng, et al. Graph convolutional networks for temporal action localization. ICCV, 2019
- Bian, Tian, et al. "Rumor Detection on Social Media with Bi-Directional Graph Convolutional Networks." AAAI 2020.
- Junchi, Yu, et al. "Graph Information Bottleneck for Subgraph Recognition.", ICLR 2021.