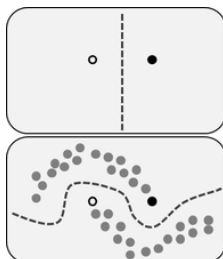




## Node Classification with Semi-supervised Learning (SSL)

- Upper: Decision boundary based on labeled data only



- Lower: Decision boundary based on labeled + unlabeled data

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

3

3

## How to support SSL for node classification?

- Approach 1.** Add a regularization term in the loss function of the classifier (e.g., **Laplacian SVM**);

$$\mathcal{L} = \mathcal{L}_0 + \mathcal{L}_{reg}, \text{ with } \mathcal{L}_{reg} = \hat{\mathbf{y}}^T \mathbf{L} \hat{\mathbf{y}} = \frac{1}{2} \sum_{i,j} A_{ij} \|\hat{\mathbf{y}}_i - \hat{\mathbf{y}}_j\|^2$$

- $\mathcal{L}_0$  is the training-set loss w.r.t. the labeled nodes in the graph;
- $\mathcal{L}_{reg}$  is the **smoothness penalty** of label propagation over unlabeled data;
- $\hat{\mathbf{y}}$  is the system-predicted scores w.r.t. the label of yes.
- Approach 2.** Use Graph Neural Networks (GNNs) for node embeddings (this lecture)
  - Geometrically move the node embeddings to be closer if they have a shared label.**
  - In other words, making the data **easily separable** for the classifiers.

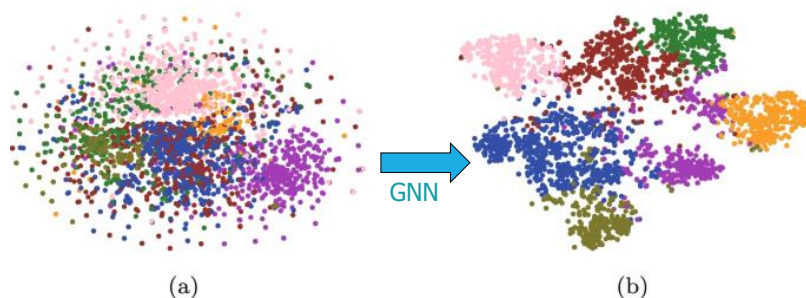
3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

4

4

## Task-aware node embeddings via GNN



3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

5

5

## Difference from the methods in previous lecture

- **Laplacian Eigenmaps and DeepWalk**
  - Both methods ignore the node-specific features and labels during embedding
- **Graph Neural Networks (GCN, GAT, etc.)**
  - Jointly leveraging node-specific features, link structures, and a training-set labels of nodes.

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

6

6

## Two types of GNNs for node classification

- GNNs with **Spatial** Convolution
  - ICLR 2018] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, **Yoshua Bengio**. "Graph attention networks" (**GAT**)
- GNNs with **Spectral** Graph Convolution
  - [ICLR 2017] Thomas N. Kipf and **Max Welling**. "Semi-supervised classification with graph convolutional networks" (**GCN**)

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

7

7

## Notation

- $G = (V, A, X)$  is an input graph (e.g., a citation graph);
  - $V$  is the set of nodes and  $|V| = n$  (e.g., each node is an article in a citation graph);
  - $A_{n \times n}$  is the adjacency matrix with non-negative elements and **zero at diagonal**;
  - $X \in \mathbb{R}^{n \times d}$  is the input matrix of **node-specific features** (one row per node, denoted as  $X_i$ ), e.g., the BERT-produced embedding per article;
- $Z \in \mathbb{R}^{n \times d'}$  is the matrix of the GNN-produced node embeddings;
- $Y_i \in \{0, 1\}^K$  is the label vector of node  $i$  for  $i = 1, \dots, n$ ;
- $\hat{Y}_i \in \mathbb{R}^K$  is the vector of system-predicted scores of labels given node  $i$ ;
- $\mathcal{D}_l = \{(X_i, Y_i)\}$  is a training set of labeled nodes.

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

8

8

## Graph Attention Networks (GAT) [P. Velickovic et al., ICLR 2018]

- A multi-head transformer module at **each layer**

- Input  $H = \{\vec{h}_1, \vec{h}_2, \dots, \vec{h}_n\}$   $\vec{h}_i \in \mathbb{R}^d$

- Output  $H' = \{\vec{h}'_1, \vec{h}'_2, \dots, \vec{h}'_n\}$   $\vec{h}'_i \in \mathbb{R}^{d'}$

- Attention between two nodes

$$e_{ij} = \theta^T [W\vec{h}_i \parallel W\vec{h}_j] \quad \text{"||" is the concatenation operator;}$$

$$\alpha_{ij} = \frac{\exp(\text{LeakReLU}(e_{ij}))}{\sum_{j' \in \mathcal{N}_i} \exp(\text{LeakReLU}(e_{ij'}))} \quad \mathcal{N}_i \text{ is the neighborhood of node } i;$$

$W \in \mathbb{R}^{d' \times d}$  and  $\theta \in \mathbb{R}^{2d'}$  are the trainable model parameters for the layer;

- Time complexity per attention assuming  $|\mathcal{N}_i| = k < d'$ :  $O(dd') + O(d') + O(k) = O(dd')$ .

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

9

9

## Multi-head Node Embedding in GAT

- Intermediate layer

$$\vec{h}'_i := \parallel_{m=1}^M \sigma \left( \sum_{j \in \mathcal{N}_i} \alpha_{ij}^m W^m \vec{h}_j \right) \quad (\vec{h}'_i \in \mathbb{R}^{Md'})$$

$\sigma(\cdot)$  is a non-linear operator (e.g., Exponentially Linear Unit or ELU).

- Final later

$$\vec{h}'_i = \sigma \left( \frac{1}{M} \sum_{m=1}^M \sum_{j \in \mathcal{N}_i} \alpha_{ij}^m W^m \vec{h}_j \right) \quad (\vec{h}'_i \in \mathbb{R}^{d'})$$

$\sigma(\cdot)$  is SoftMax or sigmoid activation for multi-class or multi-label problems.

- Model Training (for multi-class classification)

- Minimizing a cross-entropy loss over labeled training instances with Adam SGD.

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

10

10

## Evaluation Settings for Node Classification

### Transductive Learning

- Train the model with a **partially labeled graph** and predict the labels for the unlabeled nodes.
- Citation Graphs: Cora, Citeseer and PutMed (with **2.7k-19k nodes**) where the nodes are documents, the edges are citation links, and each document has one and only one label (**multi-class problems**).

### Inductive Learning

- Train the model with a set of labeled graphs, **predict the node labels on new graphs**;
- Using 20 PPI (protein-protein interaction) graphs for training, 2 PPI graphs for validation, and 2 PPI graphs for testing

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

11

11

## Node Classification Results in Transductive Learning (P. Velickovic et al., ICLR 2018)

Table 2: Summary of results in terms of **classification accuracies, for Cora, Citeseer and Pubmed**. GCN-64\* corresponds to the best GCN result computing 64 hidden features (using ReLU or ELU).

<i>Transductive</i>				
Method	Cora	Citeseer	Pubmed	
MLP	55.1%	46.5%	71.4%	MLP: ignoring graph structure
ManiReg (Belkin et al., 2006)	59.5%	60.1%	70.7%	
SemiEmb (Weston et al., 2012)	59.0%	59.6%	71.7%	
LP (Zhu et al., 2003)	68.0%	45.3%	63.0%	
DeepWalk (Perozzi et al., 2014)	67.2%	43.2%	65.3%	DeepWalk+LR(OVA): ignoring node features
ICA (Lu & Getoor, 2003)	75.1%	69.1%	73.9%	
Planetoid (Yang et al., 2016)	75.7%	64.7%	77.2%	
Chebyshev (Defferrard et al., 2016)	81.2%	69.8%	74.4%	GCN with signal passing over the given links in the graph (next)
GCN (Kipf & Welling, 2017)	81.5%	70.3%	<b>79.0%</b>	
MoNet (Monti et al., 2016)	81.7 ± 0.5%	—	78.8 ± 0.3%	
GCN-64*	81.4 ± 0.5%	70.9 ± 0.5%	<b>79.0 ± 0.3%</b>	GAT: multi-layer transformer module with graph-restricted attentions
<b>GAT (ours)</b>	<b>83.0 ± 0.7%</b>	<b>72.5 ± 0.7%</b>	<b>79.0 ± 0.3%</b>	

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

12

12

# Node Classification Results in Inductive Learning

(P. Velickovic et al., ICLR 2018)

Table 3: Summary of results in terms of **micro-averaged  $F_1$  scores, for the PPI dataset**. GraphSAGE\* corresponds to the best GraphSAGE result we were able to obtain by just modifying its architecture. Const-GAT corresponds to a model with the same architecture as GAT, but with a constant attention mechanism (assigning same importance to each neighbor; GCN-like inductive operator).

<i>Inductive</i>	
Method	PPI
Random	0.396
MLP	0.422
GraphSAGE-GCN (Hamilton et al., 2017)	0.500
GraphSAGE-mean (Hamilton et al., 2017)	0.598
GraphSAGE-LSTM (Hamilton et al., 2017)	0.612
GraphSAGE-pool (Hamilton et al., 2017)	0.600
GraphSAGE*	0.768
Const-GAT (ours)	0.934 $\pm$ 0.006
<b>GAT (ours)</b>	<b>0.973 <math>\pm</math> 0.002</b>

13

# Semi-supervised Classification with Graph Convolution Network (GCN) [Kipf and Welling, ICLR 2017]

## Key Ideas & Contributions

- Using a (two-layer) neural network with **graph convolution** for node classification
- Propagating node signals via (a modified version of) the **adjacent matrix**
- Producing **task-aware node embeddings** at the hidden layers with learnable parameters
- Motivated by the **Spectral Graph Convolution** (based on the **eigen-decomposition of the graph Laplacian**) and the 1<sup>st</sup>-order **Chebyshev approximation** to avoid explicit eigendecomposition
- Making the 1<sup>st</sup> connection of GCN with Weisfeiler-Lehman (WL) test on graph isomorphism (more discussions with GIN later)

14

## Multi-layer GCN [Kipf and Welling, ICLR 2017]

- **Input layer**

$$H^{(0)} := X = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \quad (\mathbf{x}_i \text{ is the input feature vector of node } i.)$$

- **Intermediate layers**

$$H^{(l+1)} := \sigma \left( \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)} \right) \quad \text{for } l = 0, 1, \dots, L-1$$

where  $H^{(l)} \in R^{n \times d_l}$  is the matrix of node embeddings at layer  $l$ ;

$W^{(l)} \in R^{d_l \times d_{l+1}}$  are the layer-specific matrix of learnable parameters;

$\tilde{A} = A + I_n$  is the adjacency matrix of the graph with **added self-connections** ( $I_n$ );

$\tilde{D} = \text{diag}(\tilde{D}_{ii})_{i=1, \dots, n}$  with  $\tilde{D}_{ii} = \sum_{j=1}^n \tilde{A}_{ij}$ , aggregating the degree information of nodes.

## GCN vs. GAT at an Intermediate Layer

- **GCN**  $H^{(l+1)} := \sigma \left( \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)} \right)$

- The red part does not have learnable model parameters (task-agnostic).

- **GAT**  $\vec{h}_i^{(l+1)} := \parallel_{m=1}^M \sigma \left( \sum_{j \in \mathcal{N}_i} \alpha_{ij}^{(m,l)} W^{(m,l)} \vec{h}_j^{(l)} \right)$  (slide #9)

$$H^{(l+1)} := \parallel_{m=1}^M \sigma \left( \mathcal{A}_G^{(m,l)} H^{(l)} W^{(m,l)} \right) \quad (\text{compact formula})$$

- $\mathcal{A}_G^{(m,l)}$  is the graph-masked attention matrix (based on learnable node embeddings).



## Recap of Attention in GAT (slide #8)

- Attention between two nodes

$$e_{ij} = \theta^T [W \vec{h}_i \parallel W \vec{h}_j]$$

“ $\parallel$ ” is the concatenation operator;

$$\alpha_{ij} = \frac{\exp(\text{LeakReLU}(e_{ij}))}{\sum_{j' \in \mathcal{N}_i} \exp(\text{LeakReLU}(e_{ij'}))}$$

$\mathcal{N}_i$  is the neighborhood of node  $i$ ;

$W \in R^{d \times d'}$  and  $\theta \in R^{2d'}$  are trainable model parameters.

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

17

17

## Design Choices of Signal Passing Mechanism

- Multi-layer Perceptron (MLP)  $H^{(l+1)} := \sigma(H^{(l)} W^{(l)})$  (1)

- Signal passing via  $A_{sym} = D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$   $H^{(l+1)} := \sigma(A_{sym} H^{(l)} W^{(l)})$  (2)

- Signal passing via  $\tilde{A} = I_n + A_{sym}$   $H^{(l+1)} := \sigma(\tilde{A} H^{(l)} W^{(l)})$  (3)

- “Renormalization trick”  $\hat{A} = \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$   $H^{(l+1)} := \sigma(\hat{A} H^{(l)} W^{(l)})$  (4) Empirical Winner

- Signal passing via graph Laplacian  $L_{sym} = I_n - D^{-\frac{1}{2}} A D^{-\frac{1}{2}} = U \Lambda U^T$   
 $H^{(l+1)} := \sigma(U g_{\theta}(\Lambda) U^T H^{(l)})$  (5)

- Chebyshev filter (K-th order polynomial approximation of formula 5)  
 $H^{(l+1)} := \sigma\left(\sum_{k=0}^K T_k(\tilde{L}) H^{(l)} W_k^{(l)}\right)$  for  $K = 0, 1, 2, 3$  (6)

Theoretically & Technically Motivated (for “spectral” analysis)

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

18

18

# GCN Model Training for Node Classification

(on Multi-class Datasets Cora, Citeseer, PubMed: one label per document)

- Labeled training data  $\mathcal{D}_l = \{(X_i, Y_i)\}$ 
  - $X_i \in \mathbb{R}^d$  is the feature vector of node  $i$  (a document);
  - $Y_i \in \mathbb{R}^K$  is a 0-1 valued vector for a corresponding category label.
- Output  $\hat{Y} \in \mathbb{R}^{n \times K}$  ( $n$  is the number of nodes in the graph)

$$\hat{Y} = f(X, A) = \text{softmax}(\hat{A} \text{ReLU}(\hat{A} X W^{(0)}) W^{(1)})$$

where  $\hat{A} = \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$  and  $\tilde{A} = A_{\text{sym}} + I_n$  (adjacency matrix with a re-normalization trick)

- Loss function (cross-entropy) for multi-class problem

$$\mathcal{L} = -\sum_{Y_i \in \mathcal{D}_l} \sum_{j=1}^K Y_{ij} \ln \hat{Y}_{ij}$$

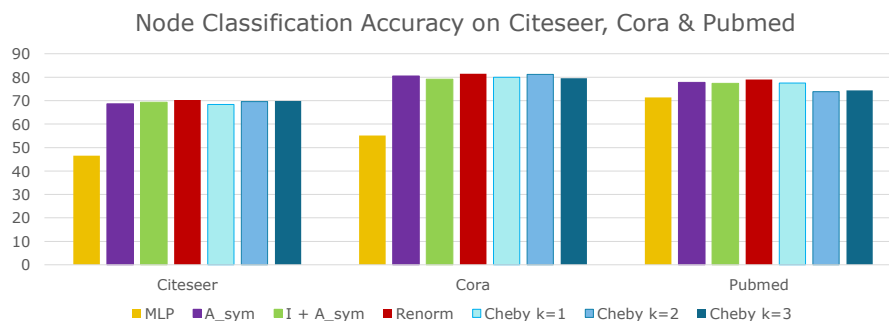
3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

19

19

## Ablation Tests [Kipf and Welling, ICLR 2017]



- Graph-based signal passing is much better than MLP (no signal passing).
- Using adjacency matrix (or modified versions) is at least as good as using Chebyshev filters.

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

20

20

## Comparative Eval Results [Kipf and Welling, ICLR 2017]

Results are summarized in Table 2. Reported numbers denote classification accuracy in percent. For ICA, we report the mean accuracy of 100 runs with random node orderings. Results for all other baseline methods are taken from the Planetoid paper (Yang et al., 2016). Planetoid\* denotes the best model for the respective dataset out of the variants presented in their paper.

Table 2: Summary of results in terms of classification accuracy (in percent).

Method	Citeseer	Cora	Pubmed	NELL
ManiReg [3]	60.1	59.5	70.7	21.8
SemiEmb [28]	59.6	59.0	71.1	26.7
LP [32]	45.3	68.0	63.0	26.5
DeepWalk + LR(One vs. All) [22]	43.2	67.2	65.3	58.1
ICA [18]	69.1	75.1	73.9	23.1
Planetoid* [29]	64.7 (26s)	75.7 (13s)	77.2 (25s)	61.9 (185s)
<b>GCN (this paper)</b>	<b>70.3 (7s)</b>	<b>81.5 (4s)</b>	<b>79.0 (38s)</b>	<b>66.0 (48s)</b>
GCN (rand. splits)	67.9 ± 0.5	80.1 ± 0.5	78.9 ± 0.7	58.4 ± 1.7

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

21

21

## Relative performance in node classification

- $GNN$  ( $GAT$  or  $GCN$ )  $\gg$   $DeepWalk$   $\gg$   $MLP$  (Why?)
- $GAT > GCN$  (Why?)

	Using node-specific features?	Using graph structure?	Link weights learnable from labeled data?
DeepWalk+LR	×	✓	
MLP	✓	×	
GCN	✓	✓	
GAT	✓	✓	

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

22

22

## Appendix p1. Graph Laplacian and its Spectrum

- Denoting signal  $x \in \mathbb{R}^n$  and filter  $g_\theta = \text{diag}(\theta)$  parameterized by  $\theta \in \mathbb{R}^n$  in the Fourier domain, the spectral convolutions on a graph is defined as

$$g_\theta * x = U g_\theta U^T x$$

where  $U$  is the matrix of eigenvectors of graph Laplacian  $L_{\text{sym}} = I_n - D^{-\frac{1}{2}} A D^{-\frac{1}{2}} = U \Lambda U^T$ ;

$\Lambda$  is the diagonal matrix of the eigenvalues (the spectrum) sorted in a descending order;

$U^T x$  is considered as being the Fourier transform of  $x$  (projecting it onto the eigenvectors);

$g_\theta = g_\theta(\Lambda)$  is a function of the eigenvalues of  $L$ , as the convolution filter in the eigen space;

$y = g_\theta(U^T x)$  is the feature map obtained by applying filter  $g_\theta$  to  $U^T x$  in the Fourier domain;

$U g_\theta U^T x = U y$  is the inverse Fourier transform of  $y$  back to the original space of  $x$ ;

$\theta \in \mathbb{R}^n$  are the trainable model parameters by a Graph Convolutional Network (GCN).

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

23

23

## Appendix p2. Computational Cost of Graph Convolution

- Directly computing spectral convolution takes  $O(n^2)$  time

$$g_\theta * x = U g_\theta U^T x$$

- Furthermore, solving the eigendecomposition of  $L$  is prohibitively expensive for large graphs.
- Remedy: approximating  $g_\theta(\Lambda)$  with Chebyshev filter of the  $K$ 'th order as

$$g_{\theta'}(\Lambda) \approx \sum_{k=0}^K \theta'_k T_k(\tilde{\Lambda}) \quad \text{with} \quad \tilde{\Lambda} = \frac{2}{\lambda_{\max}} \Lambda - I_n \quad (1)$$

- With  $K = 1, 2$  and  $3$ , the Chebyshev filters are much more efficient to compute (slide #28).

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

24

24

## Appendix p3. Chebyshev Polynomials with 1-D input

- Recursive Definition

$$T_k(v) = 2vT_{k-1}(v) - T_{k-2}(v) \text{ for } v \in R$$

$$\text{with } T_0(v) = 1; T_1(v) = v;$$

- Calculation by induction

$$T_2(v) = 2v^2 - 1;$$

$$T_3(v) = 4v^3 - 2v - v = 4v^3 - 3v$$

$$\vdots$$

- Resulted  $\sum_{k=0}^K \theta_k T_k(v)$  is a weighted combination of  $v^k$  for  $k = 0, 1, \dots, K$ .
- Similarly,  $Z_\theta(\Lambda) := \sum_{k=0}^K \theta_k T_k(\tilde{\Lambda})$  is a weighted combination of  $\tilde{\Lambda}^k$ s (diagonal matrices).

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

25

25

## Appendix p4. Chebyshev Polynomials with $\tilde{\Lambda}$ input

- Chebyshev polynomials of the  $K^{th}$  order ( $K < n$ )

$$g_{\theta'}(\Lambda) \approx \sum_{k=0}^K \theta'_k T_k(\tilde{\Lambda}) \text{ with } \tilde{\Lambda} = \frac{2}{\lambda_{max}} \Lambda - I_n$$

$$\text{with } T_0(\tilde{\Lambda}) = I_n;$$

$$T_1(\tilde{\Lambda}) = \tilde{\Lambda};$$

$$\vdots$$

$$T_k(\tilde{\Lambda}) = 2\tilde{\Lambda}T_{k-1}(\tilde{\Lambda}) - T_{k-2}(\tilde{\Lambda})$$

where  $\theta' = (\theta'_0, \theta'_1, \dots, \theta'_K)$  are the trainable Chebyshev coefficients.

3/14/2024

@Yiming Yang, 11-741 S24 GNNs Part 1

26

26

## Appendix p5. Approximated Graph Convolution with the Chebyshev Filter

$$g_\theta * x = U g_\theta U^T x \quad (2)$$

$$\approx U \sum_{k=0}^K \theta'_k T_k(\tilde{\Lambda}) U^T x = \sum_{k=0}^K \theta'_k U T_k(\tilde{\Lambda}) U^T x \quad (3)$$

$$= \sum_{k=0}^K \theta'_k T_k(\underbrace{U \tilde{\Lambda} U^T}_{\tilde{L}}) x = \sum_{k=0}^K \theta'_k T_k(\tilde{L}) x \quad (4)$$

$$\tilde{L} = U \tilde{\Lambda} U^T = U \left( \frac{2}{\lambda_{max}} \Lambda - I_n \right) U^T = \frac{2}{\lambda_{max}} L - I_n$$

- Notice that the eigendecomposition of  $L$  is no longer needed (as  $U$  and  $\Lambda$  disappeared in 4).
- $\tilde{L}$  can be highly sparse: the number of non-zero elements equals to  $|\mathcal{E}|$  (the number of edges in the graph).
- Time complexity in formula 4 is  $O(K|\mathcal{E}|)$  with sparse  $\tilde{L}$ , compared to  $O(n^2)$  in formula 2 with dense  $U$ .