GNN Architectures

AI602: Recent Advances in Deep Learning

Lecture 7

Slide made by

Insu Han and Sungsoo Ahn

KAIST EE

Algorithmic Intelligence Lab

Motivation

 Successful Deep Neural Networks (DNNs) have mainly focused on grid or sequence data, e.g., CNN and RNN.



Motivation

- Successful Deep Neural Networks (DNNs) have mainly focused on grid or sequence data, e.g., CNN and RNN.
- Many practical data are often not on the grid (i.e., non-Euclidean space)
 - E.g., distance between two nodes are not defined
- These data (e.g., molecule and social networks) can be represented by graphs



Motivation

- Practical problems of graph data
 - Community detection
 - Friend recommendation
 - Toxic detection



- \rightarrow Link prediction
- \rightarrow Graph classification



Node classification

Link prediction

Graph classification

• How can solve graph problems using DNNs?

Algorithmic Intelligence Lab

Next, Graph Neural Networks

1. Graph Neural Networks (GNNs)

- Basics
- Convolution and pooling
- Graph convolutional networks (GCNs)

2. Advanced Techniques

- Sampling neighbors for regularization
- Learning attentions of neighbors
- Extracting edge and graph features

3. Application to Combinatorial Optimization

- Boolean satisfiability
- Maximum independent set

1. Graph Neural Networks (GNNs)

- Basics
- Convolution and pooling
- Graph convolutional networks (GCNs)
- 2. Advanced Techniques
 - Sampling neighbors for regularization
 - Learning attentions of neighbors
 - Extracting edge and graph features
- 3. Application to Combinatorial Optimization
 - Boolean satisfiability
 - Maximum independent set

Recap: Deep Learning = Feature Learning

- Deep learning extracts features from the input data
 - For **image** data, deep networks finds high-level features that used for various image tasks, e.g., classification, object detection and segmentation.



- Deep learning extracts features from the input data
 - For **image** data, deep networks finds high-level features that used for various image tasks, e.g., classification, object detection and segmentation.
 - For **graph** data, we want to find **good features** for graph tasks, e.g., node classification, link prediction and graph classification.



• How can DNNs extract features on graph dataset?

- Consider (undirected) graph G := (V, E) given by adjacency matrix $A \in \mathbb{R}^{|V| \times |V|}$
- Each node feature is given by $X \in \mathbb{R}^{|V| \times d}$
 - For example, node feature can include the number of neighbors of each node

1 if node a connects to node b, 0 otherwise



Algorithmic Intelligence Lab

- Consider (undirected) graph G := (V, E) given by adjacency matrix $A \in \mathbb{R}^{|V| \times |V|}$
- Each node feature is given by $X \in \mathbb{R}^{|V| \times d}$
 - For example, node feature can include the number of neighbors of each node
- Graph Neural Networks (GNNs) can extract graph node features



- Consider convolution on image data
 - Weighted sum of features of neighbor nodes
 - The feature is multiplied with individual parameter then aggregated to next layer
 - The parameters are shared for all pixels



 $oldsymbol{x}_{e}^{(\ell+1)} = \sigma \left(oldsymbol{ heta}_{a}^{(\ell)} oldsymbol{x}_{a}^{(\ell)} + \cdots + oldsymbol{ heta}_{i}^{(\ell)} oldsymbol{x}_{i}^{(\ell)}
ight)$

individual weighted sum

- Consider **convolution** on image data
 - Weighted sum of features of neighbor nodes
 - The feature is multiplied with individual parameter then aggregated to next layer
 - The parameters are shared for all pixels



→ non-linear function

$$oldsymbol{x}_{e}^{(\ell+1)}= oldsymbol{\sigma}\left(oldsymbol{ heta}_{a}^{(\ell)}oldsymbol{x}_{a}^{(\ell)}+\dots+oldsymbol{ heta}_{i}^{(\ell)}oldsymbol{x}_{i}^{(\ell)}
ight)$$

individual weighted sum

- In graph, each node has different number of neighbor nodes
- No fixed order of neighbor nodes



- Consider convolution on image data
 - Weighted sum of features of neighbor nodes
 - The feature is multiplied with individual parameter then aggregated to next layer
 - The parameters are shared for all pixels



→ non-linear function

$$oldsymbol{x}_{e}^{(\ell+1)}= oldsymbol{\sigma}\left(oldsymbol{ heta}_{a}^{(\ell)}oldsymbol{x}_{a}^{(\ell)}+\dots+oldsymbol{ heta}_{i}^{(\ell)}oldsymbol{x}_{i}^{(\ell)}
ight)$$

individual weighted sum

• Basic approach in GNNs – share the neighbor weights



 $\boldsymbol{x}_{c}^{(\ell+1)} = \sigma \begin{pmatrix} \boldsymbol{\theta}_{0}^{(\ell)} \boldsymbol{x}_{c}^{(\ell)} + \boldsymbol{\theta}_{1}^{(\ell)} \frac{1}{|\mathcal{N}_{c}|} \sum_{j \in \mathcal{N}_{c}} \boldsymbol{x}_{j}^{(\ell)} \end{pmatrix}$ e.g., non-linear function, weight for neighbors

long short-term memory (LSTM)

Algorithmic Intelligence Lab

- [Scarselli et al., 2008] propose a GNN with a single layer (similar to RNN)
- Node embedding (feed-forward)
 - Node features is computed iteratively until converge (or maximum T steps)

$$oldsymbol{x}_{c}^{(\ell+1)} = \sigma \left(oldsymbol{ heta}_{0}oldsymbol{x}_{c}^{(\ell)} + oldsymbol{ heta}_{1}rac{1}{|\mathcal{N}_{c}|}\sum_{j\in\mathcal{N}_{c}}oldsymbol{x}_{j}^{(\ell)}
ight)$$



- [Scarselli et al., 2008] propose a GNN with a single layer (similar to RNN)
- Node embedding (feed-forward)
 - Node features is computed iteratively until converge (or maximum T steps)

$$oldsymbol{x}_{c}^{(\ell+1)} = \sigma \left(oldsymbol{ heta}_{0}oldsymbol{x}_{c}^{(\ell)} + oldsymbol{ heta}_{1}rac{1}{|\mathcal{N}_{c}|}\sum_{j\in\mathcal{N}_{c}}oldsymbol{x}_{j}^{(\ell)}
ight)$$

- Parameter learning (back-propagation)
 - Gradients are also back-propagated iteratively until converge (or maximum T steps)



- [Scarselli et al., 2008] propose a GNN with a single layer (similar to RNN)
- Forward/backward are computed iteratively until converge
- (-) In general, the states (node features) may not converge
- (-) The iteration makes computationally expensive, not scalable
- How can make GNNs computationally efficient?



Next, Graph Convolutional Networks

- Graph Convolutional Networks (GCNs) [Kipf et al., 2016a]
 - It performs fixed number of feed-forwards (e.g., 2-layer convolution)
 - Parameters for self and neighbor are tied (shared) as a single parameter
 - Normalized the features with the products of degree,





- Graph Convolutional Networks (GCNs) [Kipf et al., 2016a]
 - It performs fixed number of feed-forwards (e.g., 2-layer convolution)
 - Parameters for self and neighbor are tied (shared) as a single parameter
 - Normalized the features with the products of degree,



- (+) Number of parameters is reduced with more parameter sharing
- (+) No need to recursive operations (easy to optimize)
- (+) Easy to implement using matrix multiplication, requiring $\mathcal{O}(|E|)$ operations

→ diagonal matrix

$$X^{(\ell+1)} = \sigma \left(D^{-1/2} (A+I) D^{-1/2} X^{(\ell)} \Theta^{(\ell)} \right)$$

sparse matrix, $\mathcal{O}(|E| + |V|)$ non-zero entries

- Graph Convolutional Networks (GCNs) [Kipf et al., 2016a]
 - It performs fixed number of feed-forwards (e.g., 2-layer convolution)
 - Parameters for self and neighbor are tied (shared) as a single parameter
 - Normalized the features with the products of degree,



- (+) Number of parameters is reduced with more parameter sharing
- (+) No need to recursive operations (easy to optimize)
- (+) Easy to implement using matrix multiplication, requiring $\mathcal{O}(|E|)$ operations
- (+) Residual connections can improve the performance

$$X^{(\ell+1)} = \sigma \left(D^{-1/2} (A+I) D^{-1/2} X^{(\ell)} \Theta^{(\ell)} \right) + X^{(\ell)}$$

residual connection

Algorithmic Intelligence Lab

- Overview of CNNs, GNNs and GCNs
 - CNNs have individual parameters for all neighbors and self nodes
 - GNNs have 2 type of parameters; only parameters for neighbor nodes are shared
 - GCNs have a single parameter shared for both self and neighbor nodes



- Graph Convolutional Networks (GCNs) [Kipf et al., 2016a]
- Goal classify the paper category (e.g., stat.ML, cs.LG, and etc) for citation networks (graph node classification)
- Model 2 layer GCN $Z = \operatorname{softmax} \left(\widehat{A} \operatorname{ReLU} \left(\widehat{A} X \Theta^{(0)} \right) \Theta^{(1)} \right)$

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 [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)
GCN (this paper)	70.3 (7s)	81.5 (4s)	79.0 (38s)	66.0 (48s)

Test accuracy (training time) for graph node classification

- Graph Convolutional Networks (GCNs) [Kipf et al., 2016a]
- Goal classify the paper category (e.g., stat.ML, cs.LG, and etc) for citation networks (graph node classification)
- Model 2 layer GCN $Z = \operatorname{softmax} \left(\widehat{A} \operatorname{ReLU} \left(\widehat{A} X \Theta^{(0)} \right) \Theta^{(1)} \right)$
- Why 2 layer is enough?
 - Most of graph dataset can be covered with 4 hop-neighbors
 - Large number of layers can cause over-smoothing on the node features



- Graph Convolutional Networks (GCNs) [Kipf et al., 2016a]
 - (+) Number of parameters is reduced with more parameter sharing
 - (+) No need to recursive operations (easy to optimize)
 - (+) Easy to implement using matrix multiplication, requiring $\mathcal{O}(|E|)$ operations
 - (+) Empirically works well even for a small number of layers
- Some limitations on GCNs
 - (-) Can be overfit easily
 - (-) Still expensive for large-scale graphs
 - (-) Cannot embedding edge features
- How can overcome these drawbacks?

Table of Contents

1. Graph Neural Networks (GNNs)

- Basics
- Convolution and pooling
- Graph convolutional networks (GCNs)

2. Advanced Techniques

- Sampling neighbors for regularization
- Learning attentions of neighbors
- Extracting edge and graph features

3. Application to Combinatorial Optimization

- Boolean satisfiability
- Maximum independent set

- **Problem** GCNs can be overfit easily
- **Solution** Train the model stochastically based on sampling neighbors

- Problem GCNs can be overfit easily
- Solution Train the model stochastically based on sampling neighbors
- Graph Sampling and Aggregation (GraphSAGE) [Hamilton et al., 2017]
- Sample neighbor nodes randomly
 - (+) Avoid overfitting
 - (+) Training time can be reduced

Aggregate
$$\left(\left\{ \boldsymbol{x}_{j}^{(\ell)}; j \in S_{c}, S_{c} \subseteq \mathcal{N}_{c} \right\} \right)$$

neighbor sampling



- **Problem** GCNs can be overfit easily
- Solution Train the model stochastically based on sampling neighbors
- Graph Sampling and Aggregation (GraphSAGE) [Hamilton et al., 2017]
- Sample neighbor nodes randomly
 - (+) Avoid overfitting
 - (+) Training time can be reduced
- Instead of summing the self and neighbor features, concatenate them

$$\boldsymbol{x}_{c}^{(\ell+1)} = \sigma\left(\left[\boldsymbol{\theta}_{0}^{(\ell)}\boldsymbol{x}_{c}^{(\ell)}, \text{ Aggregate}\left(\left\{\boldsymbol{x}_{j}^{(\ell)}; j \in S_{c}, S_{c} \subseteq \mathcal{N}_{c}\right\}\right)\right]\right)$$



- **Problem** GCNs can be overfit easily
- Solution Train the model stochastically based on sampling neighbors
- Graph Sampling and Aggregation (GraphSAGE) [Hamilton et al., 2017]
- Sample neighbor nodes randomly
 - (+) Avoid overfitting
 - (+) Training time can be reduced
- Instead of summing the self and neighbor features, concatenate them

$$\boldsymbol{x}_{c}^{(\ell+1)} = \sigma\left(\left[\boldsymbol{\theta}_{0}^{(\ell)}\boldsymbol{x}_{c}^{(\ell)}, \text{ Aggregate}\left(\left\{\boldsymbol{x}_{j}^{(\ell)}; j \in S_{c}, S_{c} \subseteq \mathcal{N}_{c}\right\}\right)\right]\right)$$

neighbor sampling

• Aggregation function can be generalized including

Mean	$igg rac{1}{ \mathcal{N}_c } \sum_{j \in \mathcal{N}_c} oldsymbol{x}_j^{(\ell)}$	
Pool	$\max\left(\left\{W^{(\ell)}oldsymbol{x}_{j}^{(\ell)}; j\in\mathcal{N}_{c} ight\} ight)$	$W^{(\ell)}$: learnable parameter
LSTM	$\operatorname{LSTM}\left(\left[oldsymbol{x}_{j}^{(\ell)}; j\in\pi(\mathcal{N}_{c}) ight] ight)$	π : random permutation

- **Problem** GCNs can be overfit easily
- Solution Train the model stochastically based on sampling neighbors
- Graph Sampling and Aggregation (GraphSAGE) [Hamilton et al., 2017]
- Aggregation function can be {mean, pool, LSTM}
- Experiments for node classification
 - Performs better than GCNs

	Citation		Rede	lit	PPI		
Name	Unsup. F1	Sup. F1	Unsup. F1	Sup. F1	Unsup. F1	Sup. F1	
Random	0.206	0.206	0.043	0.042	0.396	0.396	
Raw features	0.575	0.575	0.585	0.585	0.422	0.422	
DeepWalk	0.565	0.565	0.324	0.324			
DeepWalk + features	0.701	0.701	0.691	0.691	—	—	
GraphSAGE-GCN	0.742	0.772	0.908	0.930	0.465	0.500	
GraphSAGE-mean	0.778	0.820	0.897	0.950	0.486	0.598	
GraphSAGE-LSTM	0.788	0.832	0.907	0.954	0.482	0.612	
GraphSAGE-pool	0.798	0.839	0.892	0.948	0.502	0.600	
% gain over feat.	39%	46%	55%	63%	19%	45%	

Next, GCNs with Edge Attention

- **Problem** GCNs use the fixed weight (or importance) of the edges
- **Solution** Use the learnable edge attention weights

$$\alpha_{ij} = \mathbf{f} \left(\mathbf{x}_i, \{ \mathbf{x}_j : j \in \mathcal{N}_i \} \right) \qquad \mathbf{x}_i^{(\ell+1)} = \sigma \left(\mathbf{\theta}^{(\ell)} \sum_{j \in \mathcal{N}_i \cup \{i\}} \alpha_{ij} \mathbf{x}_j^{(\ell)} \right)$$

- **Problem** GCNs use the fixed weight (or importance) of the edges
- **Solution** Use the learnable edge attention weights

$$\alpha_{ij} = \mathbf{f} \left(x_i, \{ x_j : j \in \mathcal{N}_i \} \right) \qquad \mathbf{x}_i^{(\ell+1)} = \sigma \left(\mathbf{\theta}^{(\ell)} \sum_{j \in \mathcal{N}_i \cup \{i\}} \alpha_{ij} \mathbf{x}_j^{(\ell)} \right)$$

• GCNs have constant edge weights, i.e.,

$$\alpha_{ij}^{(\text{GCN})} = \frac{1}{\sqrt{|\mathcal{N}_i||\mathcal{N}_j|}}$$

• GraphSAGE has binary random edge weights, i.e.,

$$\alpha_{ij}^{(\text{SAGE})} = \begin{cases} 1 & \text{w.p. } \frac{1}{|\mathcal{N}_i|} \\ 0 & \text{otherwise} \end{cases}$$

• Edge weights can be a trainable neural network

Learning Attentions of Neighbors

- Problem GCNs use the fixed weight (or importance) of the edges
- **Solution** Use the learnable edge attention weights
- Graph Attention Networks (GAT) [Veličković et al., 2017]

$$oldsymbol{x}_i^{(\ell+1)} = \sigma \left(oldsymbol{ heta}^{(\ell)} \sum_{j \in \mathcal{N}_i \cup \{i\}} lpha_{ij} oldsymbol{x}_j^{(\ell)}
ight)$$

$$\alpha_{ij} = \frac{\exp\left(\text{LeakyReLU}\left(\frac{\text{NN}([\boldsymbol{x}_i, \boldsymbol{x}_j])\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\frac{\text{NN}([\boldsymbol{x}_i, \boldsymbol{x}_k])\right)\right)}{\text{single layer neural network}}$$



- Problem GCNs use the fixed weight (or importance) of the edges
- **Solution** Use the learnable edge attention weights
- Graph Attention Networks (GAT) [Veličković et al., 2017]

$$\boldsymbol{x}_{i}^{(\ell+1)} = \sigma \left(\boldsymbol{\theta}^{(\ell)} \sum_{j \in \mathcal{N}_{i} \cup \{i\}} \alpha_{ij} \boldsymbol{x}_{j}^{(\ell)} \right)$$

$$\alpha_{ij} = \frac{\exp\left(\text{LeakyReLU}\left(\text{NN}([\boldsymbol{x}_i, \boldsymbol{x}_j])\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\text{NN}([\boldsymbol{x}_i, \boldsymbol{x}_k])\right)\right)}_{\substack{\text{single layer neural network}}}$$

Multiple attentions can be used with averaging or concatenating

$$\boldsymbol{x}_{i}^{(\ell+1)} = \sigma \left(\frac{1}{K} \sum_{k=1}^{K} \left(\boldsymbol{\theta}^{(\ell,K)} \sum_{j \in \mathcal{N}_{i} \cup \{i\}} \alpha_{ij}^{(k)} \boldsymbol{x}_{j}^{(\ell)} \right) \right)$$

concatenate/average of K independent attentions



Algorithmic Intelligence Lab

- Problem GCNs use the fixed weight (or importance) of the edges
- **Solution** Use the learnable edge attention weights
- Graph Attention Networks (GAT) [Veličković et al., 2017]

$$oldsymbol{x}_i^{(\ell+1)} = \sigma \left(oldsymbol{ heta}^{(\ell)} \sum_{j \in \mathcal{N}_i \cup \{i\}} lpha_{ij} oldsymbol{x}_j^{(\ell)}
ight)$$

$$\alpha_{ij} = \frac{\exp\left(\text{LeakyReLU}\left(\frac{\text{NN}([\boldsymbol{x}_i, \boldsymbol{x}_j])\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\frac{\text{NN}([\boldsymbol{x}_i, \boldsymbol{x}_k])\right)\right)}{\text{Single layer neural network}}$$

- Multiple attentions can be used with averaging or concatenating
- (+) Edge features can be learned
- (+) Generalize the neighbor sampling (soft version)
- (-) Traning is slower than GCN and difficult to optimize

- **Problem** GCNs use the fixed weight (or importance) of the edges
- **Solution** Use the learnable edge attention weights
- Graph Attention Networks (GAT) [Veličković et al., 2017]
- Experiments for node classification
 - Performs better than GCNs

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 ± 0.006
GAT (ours)	$\textbf{0.973} \pm 0.002$

- GNNs basically learns a good node embeddings
- How to use the node embeddings of GNNs for other tasks?
- Link (or edge) prediction
 - Edge feature can be the dot-product of adjacent node features [Kipf et al., 2016b]
 - Loss function can be the cross-entropy for binary classification



- GNNs basically learns a good node embeddings
- How to use the node embeddings of GNNs for other tasks?
- Graph-level classification
 - Graph features can be the summation or max pooling of all node features [Duvenaud et al., 2015]
 - A feed forward network can be used as a classifier
 - Graph pooling can be important to compress the graph feature



- GNNs basically learns a good node embeddings
- How to use the node embeddings of GNNs for other tasks?
- Graph-level classification
 - Graph features can be the summation or max pooling of all node features [Duvenaud et al., 2015]
 - A feed forward network can be used as a classifier
 - Graph pooling can be important to compress the graph feature
- How to graph pooling ?
 - Pooling methods can be learned with parameters [Ying et al., 2018]



Algorithmic Intelligence Lab

*source: https://arxiv.org/pdf/1806.08804.pdf 38

Differentiable pooling [Ying et al., 2018] – graph pooling can be represented using matrix multiplication on adjacency matrix and node feature matrix

$$\hat{A} = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \longrightarrow \hat{A} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$$

$$\hat{A} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$$

$$\hat{A} = S^{\top} A S$$

$$\hat{A} = S^{\top} A S$$

$$\hat{S} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}$$

- Assignment matrix has the same dimension of node features
- **Idea** learning the assignmet matrix using GNNs

 $S = \operatorname{softmax}(\operatorname{GNN}(A, X))$

• Soft pooling operation can be differentiable and shows improvement for graph classification

0/

Table of Contents

1. Graph Neural Networks (GNNs)

- Basics
- Convolution and pooling
- Graph convolutional networks (GCNs)

2. Advanced Techniques

- Sampling neighbors for regularization
- Learning attentions of neighbors
- Extracting edge and graph features

3. Application to Combinatorial Optimization

- Boolean satisfiability
- Maximum independent set

 Boolean satisfiability (SAT) problem determines the existence of interpretation satisfying the Boolean formula:

 $\begin{array}{c} \textit{Literal} & \textit{Clause} \\ \uparrow & \\ (x_1 \lor x_2 \lor x_3) \land (\bar{x}_1 \lor \bar{x}_2 \lor \bar{x}_3) \land (\bar{x}_1 \lor x_2 \lor x_3) \end{array}$

Example of boolean formula

• The example formula is satisfiable by following interpretation:

 $x_1 = \text{True}, \quad x_2 = \text{False}, \quad x_3 = \text{True}$

Example of satisfying interpretation

- It is very hard to solve in general, i.e., NP-hard
- Various applications, e.g., circuit design and automatic theorem proving
- A GNN can be trained to determine the existence [Selsam et al., 2018]

- **Step 1**: Represent the Boolean formula as a graph
 - Nodes represent literals and clauses in the formula
 - Edges are added for each literal-clause pair
 - Edges are added for literals with opposite signs



- Step 2: Train the GNN to predict satisfiability of formula (True / False)
 - This becomes a binary graph classification problem
 - Labels of graphs, i.e., true satisfiability, are obtained by existing solvers

• More details: the GNN consists of two phases



```
c_1 \leftarrow \text{LSTM}(x_1, x_2)
```

"clause-update" phase



$x_1, x_2 \leftarrow \text{LSTM}^{(1)}(c_1), \text{LSTM}^{(2)}(c_1)$

"literal-update" phase

- Long short-term memory (LSTM) can be used for embedding features
- Two phases (or layers) are repeated multiple times
- **Result:** the trained GNN is better than random guessing
 - Trained on SAT problems with up to 40 variables
 - Achieves 85% test accuracy (existing solvers get 100% correct)

• Independent set is a set of nodes, where no two nodes are adjacent



Examples of independent set (with size 3, 1, 3)

- Maximum independent set (MIS) problem search for the independent set with maximum cardinality
 - Again, hard to solve, i.e., NP-hard
 - Applications to computer vision, communication, physics, ...
- A GCN can be trained to find the MIS [Li et al., 2018]

- The naïve approach: train a GCN to predict the MIS
 - Becomes a node classification problem



Label "1" means "included in the independent set"

$$\mathcal{L} := \sum_{G} \text{cross-entropy} \left(\text{GCN}(X_G) \right) \longleftarrow \text{ Labels are given by existing solvers}$$

- **Problem:** at generation, the output sometimes have a pair of adjacent nodes
 - This violates the assumption of independent set

Maximum Independent Set Problem

- **Problem:** the output sometimes can have a pair of adjacent nodes
- Solution: label them one-by-one
- Repeatedly update the labels as follows:
 - 1. Assign **1** to the node with the highest probability of being chosen
 - 2. Assign **0** to the adjacent nodes





• Result: the trained GNN successfully solves the MIS problems

	MIS			MVC				
Name	Classic	S2V-DQN	ReduMIS	Ours	Classic	S2V-DQN	ReduMIS	Ours
ego-Facebook	993	1,020	1,046	1,046	3,046	3,019	2,993	2,993
ego-Gplus	56,866	56,603	57,394	57,394	50,748	51,011	50,220	50,220
ego-Twitter	36,235	36,275	36,843	36,843	45,071	45,031	44,463	44,463
soc-Epinions1	53,457	53,089	53,599	53,599	22,422	22,790	22,280	22,280
soc-Slashdot0811	53,009	52,719	53,314	53,314	24,351	24,641	24,046	24,046
soc-Slashdot0922	56,087	55,506	56,398	56,398	26,081	26,662	25,770	25,770
wiki-Vote	4,730	4,779	4,866	4,866	2,385	2,336	2,249	2,249
wiki-RfA	8,019	7,956	8,131	8,131	2,816	2,879	2,704	2,704
bitcoin-otc	4,330	4,334	4,346	4,346	1,551	1,547	1,535	1,535
bitcoin-alpha	2,703	2,705	2,718	2,718	1,080	1,078	1,065	1,065

Objective of MIS problems in large graphs (bold is best)

• The GNN solver performs well even when training and test graphs are different

- GNNs have similar operations with CNNs, but specialized for graph-structured dataset, i.e., neighbor aggregation
- GNNs basically extracts node features that suitable for graph tasks, e.g., node classification; From these node features, it can be used for link prediction and graph-level classicification
- GNNs are first proposed using iterative feed-forwards, after that GCNs are developed which are more efficient and perform better
- Several variants of GCNs, e.g., GraphSAGE, GAT, have been proposed to overcome drawbacks of GCNs
- GCNs practically are used for many graph problems including node classification, link prediction, graph classification and combinatorial optimization

[Scarselli et al., 2008] Scarselli, Franco, et al. "The graph neural network model." IEEE Transactions on Neural Networks. 2008.

link: https://ieeexplore.ieee.org/abstract/document/4700287

[Duvenaud et al., 2015] Duvenaud, David K., et al. "Convolutional networks on graphs for learning molecular fingerprints." Advances in neural information processing systems. 2015. link: http://papers.nips.cc/paper/5954-convolutional-networks-on-graphs-for-learning-molecular-fingerprints.pdf

[Kipf et al., 2016a] Kipf, Thomas N., and Max Welling. "Semi-supervised classification with graph convolutional networks." arXiv preprint arXiv:1609.02907. 2016. link: https://openreview.net/pdf?id=SJU4ayYgl

[Kipf et al., 2016b] Kipf, Thomas N., and Max Welling. "Variational graph auto-encoders." arXiv preprint arXiv:1611.07308. 2016.

link: https://arxiv.org/pdf/1611.07308.pdf

[Hamilton et al., 2017] Hamilton, Will, Zhitao Ying, and Jure Leskovec. "Inductive representation learning on large graphs." Advances in Neural Information Processing Systems. 2017. link: <u>https://papers.nips.cc/paper/6703-inductive-representation-learning-on-large-graphs.pdf</u>

[Veličković et alk ., 2017] Veličković, Petar, et al. "Graph attention networks." arXiv preprint arXiv:1710.10903. 2017. link: <u>https://openreview.net/forum?id=rJXMpikCZ</u>

[Ying et al., 2018] Ying, Zhitao, et al. "Hierarchical graph representation learning with differentiable pooling." Advances in Neural Information Processing Systems. 2018.

link: <u>http://papers.nips.cc/paper/7729-hierarchical-graph-representation-learning-with-differentiable-pooling.pdf</u>

[Selsam et al., 2018] Selsam et al. "Learning a SAT solver from single-bit super-vision" International Conference of Learning Representations (2018).

link: <u>https://openreview.net/pdf?id=HJMC_iA5tm</u>

[Li et al., 2018] Li et al. "Combinatorial Optimization with Graph Convolutional Networks and Guided Tree Search." Advances in Neural Information Porcessing Systems. 2018. link: <u>https://arxiv.org/pdf/1810.10659.pdf</u>