# Graph Neural Networks:
# A Review of Methods and Applications

Jie Zhou*, Ganqu Cui*, Zhengyan Zhang*, Cheng Yang, Zhiyuan Liu,
Maosong Sun
Tsinghua University
arxiv 2019

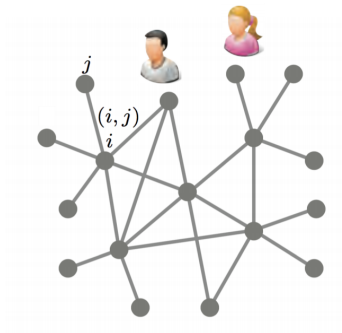https://qdata.github.io/deep2Read

Presenter : Jack Lanchantin

# Outline

# Graphs

Type of data structure which models a set of objects (vertices) and their relationships (edges).
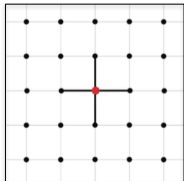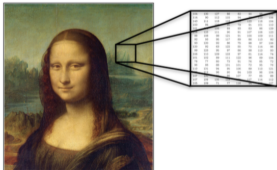
# Graphs

Type of data structure which models a set of objects (vertices) and their relationships (edges).

- **Graph** $\qquad \mathcal{G} = (\mathcal{V}, \mathcal{E})$

- **Vertices** $\qquad \mathcal{V} = \{1, \ldots, n\}$

- **Edges** $\qquad \mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$

- **Vertex weights** $\quad b_i > 0$ for $i \in \mathcal{V}$

- **Edge weights** $\qquad a_{ij} \geq 0$ for $(i, j) \in \mathcal{E}$
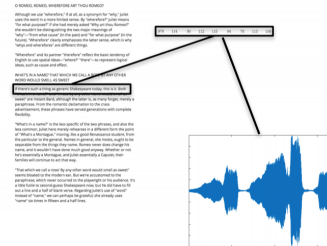
- Image, volume, video: 2D, 3D, 2D+1 Euclidean domains





2D grids

- Sentence, word, sound: 1D Euclidean domain





1D grid
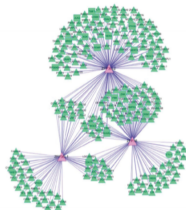
- These domains have nice regular spatial structures.
  ⇒ All ConvNet operations are math well defined and fast (convolution, pooling).

Social networks

Regulatory networks

Functional networks

3D shapes

**=**

Graphs/
Networks

- Also chemistry, NLP, physics, social science, communication networks, etc.

- Assumption: Non-Euclidean data are locally stationary and manifest hierarchical structures.
- But, how to define compositionality on graphs? (convolution and pooling on graphs)

# Outline

Naive approach: a per-node classifier

- Represent each node $v$ as vector $\mathbf{h}_v \in \mathbb{R}^s$
- Completely drop the graph structure, and classify each node individually, with a shared deep neural network classifier:

$$\mathbf{o}_v = f(\mathbf{h}_v; \mathbf{W}) \tag{1}$$

# Learning on Graphs

Naive approach: a per-node classifier

- Represent each node $v$ as vector $\mathbf{h}_v \in \mathbb{R}^s$
- Completely drop the graph structure, and classify each node individually, with a shared deep neural network classifier:

$$\mathbf{o}_v = f(\mathbf{h}_v; \mathbf{W}) \tag{1}$$



- Methods like DeepWalk also classify each node independently, but inject graph structure indirectly using learned embeddings

- Most deep learning is done this way, even if there are relationships between training examples

- Originally introduced by Scarselli et. al. (2009)
- The goal of GNNs is to learn a state embedding $\mathbf{h}_v \in \mathbb{R}^s$ which contains information of the neighborhood for each node
- $\mathbf{h}_v$ can be used to produce an output $\mathbf{o}_v$ such as the node label

- Local transition function $f$: updates each node state according to its neighborhood (shared among all nodes)
- Output function $g$: describes how the output is produced

# Graph Neural Networks

- Local transition function $f$: updates each node state according to its neighborhood (shared among all nodes)
- Output function $g$: describes how the output is produced
- Then, $\mathbf{h}_v$ and $\mathbf{o}_v$ are defined as follows:

$$\mathbf{h}_v = f(\mathbf{x}_v, \mathbf{x}_{e(v)}, \mathbf{h}_{n(v)}, \mathbf{x}_{n(v)}) \qquad (2)$$

$$\mathbf{o}_v = g(\mathbf{h}_v, \mathbf{x}_v) \qquad (3)$$

where $\mathbf{x}_v$ are the features of $v$, $\mathbf{x}_{e(v)}$ the features of its edges, $\mathbf{h}_{n(v)}$ the states of the nodes in the neighborhood of $v$, and $\mathbf{x}_{n(v)}$ the features of the nodes in the neighborhood of $v$

GNNs use the following iterative update to compute the state:

$$\mathbf{H}^{t+1} = f(\mathbf{H}^t, \mathbf{X}) \tag{4}$$

where $\mathbf{H}^t$ denotes the $t$-th iteration of $\mathbf{H}$.

# Graph Neural Network Training

- Given a GNN framework, learn the parameters of $f$ and $g$
- With the target information ($\mathbf{t}_v$ for a specific node) for the supervision, the loss can be written as:

$$loss = \sum_{i=1}^{p} (\mathbf{t}_i - \mathbf{o}_i) \tag{5}$$

  where $p$ is the number of supervised nodes.
- Weights $\mathbf{W}$ of $f$ and $g$ are updated via gradient descent

- Original GNN constrains $f$ to be a contractive map
  - Contractive map, on a metric space $(M, d)$ is a function $f$ from $M$ to itself, with the property that there is some real number $0 \leq k < 1$ such that for all $x$ and $y$ in M, $d(f(x), f(y)) \leq k \, d(x, y)$

- Original GNN constrains $f$ to be a contractive map
  - Contractive map, on a metric space $(M, d)$ is a function $f$ from $M$ to itself, with the property that there is some real number $0 \leq k < 1$ such that for all $x$ and $y$ in M, $d(f(x), f(y)) \leq k\, d(x, y)$
- This implies that the $h_v$ vectors will always converge to a unique fixed point (very restrictive)
- Impossible to inject problem-specific information into $h_v^0$ (will always converge to same value regardless of initialization)

Let $\mathbf{H}$, $\mathbf{O}$, $\mathbf{X}$, and $\mathbf{X}_N$ be the vectors constructed by stacking all the states, all the outputs, all the features, and all the node features, respectively. Then we have a compact form as:

$$\mathbf{H} = f(\mathbf{H}, \mathbf{X}) \tag{6}$$

$$\mathbf{O} = g(\mathbf{H}, \mathbf{X}_N) \tag{7}$$

# Outline

# Directed Graphs

- e.g. Knowledge graphs where there is a parent and child.
- Can use different weights for each edge type

# Heterogeneous Graphs

- e.g. Multi-model biological network data
- Can group neighbors based on types or distance

# Graphs with Edge Information

- e.g. Drug interaction data with different types of interactions (edges). Each edge also has its information like the weight or the type of the edge.
- Can use different types of weight matrices

# Graph Types

# Outline

- **Node-level outputs**: prediction for each node in the graph (e.g. classify each person in a social network)
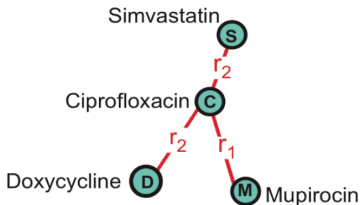- **Edge-level outputs**: prediction for each edge in the graph (e.g. classify each friendship in a social network)
- **Graph-level outputs**; prediction on the graph as a whole (e.g. classify an entire group in a social network)

- **Supervised learning for node, edge, or graph level classification**: Given full labeled networks, predict target objective
- **Semi-supervised learning for node or edge level classification**: Given a single network with partial nodes being labeled and others remaining unlabeled
- **Unsupervised learning for graph embedding**: When no class labels are available in graphs, we can learn the graph embedding in a purely unsupervised way

# Outline

- The propagation (update) step and output step are the crucial components to obtain the hidden states of nodes (or edges).

$$\mathbf{H} = f(\mathbf{H}, \mathbf{X})$$

- Variants utilize different aggregators to gather information from each node's neighbors and updaters to update nodes' hidden states

# Outline

Advances in this direction are often categorized into either spectral or spatial approaches

# Spectral Convolution Methods

- Define a convolutional operation by operating on the graph in the spectral domain, leveraging the convolution theorem
- These approaches utilise the graph Laplacian matrix, L, defined as $\mathbf{L} = \mathbf{D} - \mathbf{A}$, where $\mathbf{D}$ is the degree matrix (diagonal matrix with $\mathbf{D}_{ii} = \deg(i)$) and $\mathbf{A}$ is the adjacency matrix.

| Labeled graph | Degree matrix | Adjacency matrix | Laplacian matrix |
|---|---|---|---|
|  | $\begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$ | $\begin{pmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$ | $\begin{pmatrix} 2 & -1 & 0 & 0 & -1 & 0 \\ -1 & 3 & -1 & 0 & -1 & 0 \\ 0 & -1 & 2 & -1 & 0 & 0 \\ 0 & 0 & -1 & 3 & -1 & -1 \\ -1 & -1 & 0 & -1 & 3 & 0 \\ 0 & 0 & 0 & -1 & 0 & 1 \end{pmatrix}$ |

- Poor generalization to new/different graphs
- Graphs with variable size: spectral techniques work with fixed size graphs.
- Directed graphs: definition of directed graph Laplacian is unclear.

# Transductive Learning
(from Petar Velickovic)



Training algorithm sees all features (including test nodes)

- Algorithm does not have access to all nodes upfront. Two variations:
  1. Test nodes are (incrementally) inserted into training graphs
  2. Test graphs are disjoint and completely unseen

- Requires generalizing across arbitrary graph structures, thus many transductive methods are useless

Euclidean
space/grid

Non-Euclidean
space/graph

Fixed domain

$\Rightarrow$

Change one single
edge

Variable domain

Standard
ConvNets

Spectral graph
ConvNets

Can we still use
spectral graph
ConvNets?

Spectral NNs offer rich
families of spectral filters

$\Rightarrow$

Are spectral filters
transferable?

# Spatial Convolution Methods

- Spatial approaches define convolutions directly on the graph, operating on spatially close neighbors
- Major challenge: defining convolution with different sized neighborhoods and maintaining the local invariance of CNNs

- Different weight matrices are used for nodes with different degrees,

$$\mathbf{x} = \mathbf{h}_v + \sum_{i=1}^{\mathcal{N}_v} \mathbf{h}_i$$
$$\mathbf{h}'_v = \sigma\left(\mathbf{x}\mathbf{W}_L^{\mathcal{N}_v}\right)$$

(8)

where $\mathbf{W}_L^{\mathcal{N}_v}$ is the weight matrix for nodes with degree $\mathcal{N}_v$ at layer $L$

- **Problem**: doesn't scale to graphs with very wide degree distributions

- Transition matrices are used to define the neighborhood for nodes in DCNN. For node classification, it has

$$\mathbf{H} = f\left(\mathbf{W}^c \odot \mathbf{P}^*\mathbf{X}\right) \tag{9}$$

- $\mathbf{X}$ is an $N \times F$ matrix of input features ($N$ is the number of nodes and $F$ is the number of features)
- $\mathbf{P}^*$ is an $N \times K \times N$ tensor which contains the power series $\{\mathbf{P}, \mathbf{P}^2, ..., \mathbf{P}^K\}$ of matrix $\mathbf{P}$ (where $\mathbf{P}$ is the degree- normalized transition matrix from the adjacency matrix $\mathbf{A}$)
    - $\mathbf{P}^*$ transforms each node into a $K \times F$ diffusion convolutional representation

- General inductive framework
- Restricts every degree to be the same (by sampling a fixed-size set a node's local neighborhood, during both training and inference)

$$\mathbf{h}^t_{\mathcal{N}_v} = \mathrm{AGGREGATE}_t \left( \{ \mathbf{h}^{t-1}_u, \forall u \in \mathcal{N}_v \} \right)$$
$$\mathbf{h}^t_v = \sigma \left( \mathbf{W}^t \cdot [\mathbf{h}^{t-1}_v \| \mathbf{h}^t_{\mathcal{N}_v}] \right)$$

(10)



1.Sample neighborhood    2. Aggregate feature information from neighbors    3. Predict node labels

- General inductive framework
- Restricts every degree to be the same (by sampling a fixed-size set a node's local neighborhood, during both training and inference)

$$\mathbf{h}_{\mathcal{N}_v}^t = \text{AGGREGATE}_t \left( \{ \mathbf{h}_u^{t-1}, \forall u \in \mathcal{N}_v \} \right)$$
$$\mathbf{h}_v^t = \sigma \left( \mathbf{W}^t \cdot [\mathbf{h}_v^{t-1} \| \mathbf{h}_{\mathcal{N}_v}^t] \right)$$

(10)



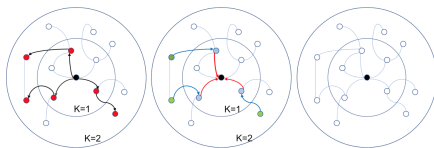1.Sample neighborhood    2. Aggregate feature information    3. Predict node labels
                         from neighbors

- Key is that $\mathbf{W}^t$ is not shared across time steps $T$

# Spatial Methods: Summary

| Variant | Aggregator | Updater |
|---------|-----------|---------|
| Convolutional networks in [33] | $\mathbf{h}_{\mathcal{N}_v}^t = \mathbf{h}_v^{t-1} + \sum_{k=1}^{\mathcal{N}_v} \mathbf{h}_k^{t-1}$ | $\mathbf{h}_v^t = \sigma(\mathbf{h}_{\mathcal{N}_v}^t \mathbf{W}_L^{\mathcal{N}_v})$ |
| DCNN | Node classification: <br> $\mathbf{N} = \mathbf{P}^* \mathbf{X}$ <br><br> Graph classification: <br> $\mathbf{N} = 1_N^T \mathbf{P}^* \mathbf{X}/N$ | $\mathbf{H} = f\left(\mathbf{W}^c \odot \mathbf{N}\right)$ |
| GraphSAGE | $\mathbf{h}_{\mathcal{N}_v}^t = \mathrm{AGGREGATE}_t\left(\{\mathbf{h}_u^{t-1}, \forall u \in \mathcal{N}_v\}\right)$ | $\mathbf{h}_v^t = \sigma\left(\mathbf{W}^t \cdot [\mathbf{h}_v^{t-1} \| \mathbf{h}_{\mathcal{N}_v}^t]\right)$ |

# Outline

- Extension of original GNN (Scarselli et. al. 2009)
- Propagate for $T$ steps, but do not restrict the propagation model to be contractive
  - Use gating in the propagation step to alleviate gradient issues

# Gating

Initial idea: update $\mathbf{h}_v^t$ using neighborhood summation and tanh

$$\mathbf{a}_v^t = \mathbf{b} + \sum_{j \in \mathcal{N}_v} \mathbf{h}_j^{t-1} \tag{11}$$

$$\mathbf{h}_v^t = \tanh(\mathbf{W}\mathbf{a}_v^t) \tag{12}$$

Now, extend this to incorporate gating mechanisms, to prevent full overwrite of $\mathbf{h}_v^{t-1}$ by $\mathbf{h}_v^t$

Gating mechanism defined similar to LSTM:

$$\mathbf{a}_v^t = \mathbf{b} + \sum_{j \in \mathcal{N}_v} \mathbf{h}_j^{t-1} \tag{13}$$

$$\mathbf{z}_v^t = \sigma \left( \mathbf{W}^z \mathbf{a}_v^t + \mathbf{U}^z \mathbf{h}_v^{t-1} \right) \tag{14}$$

$$\mathbf{r}_v^t = \sigma \left( \mathbf{W}^r \mathbf{a}_v^t + \mathbf{U}^r \mathbf{h}_v^{t-1} \right) \tag{15}$$

$$\widetilde{\mathbf{h}_v^t} = \tanh \left( \mathbf{W} \mathbf{a}_v^t + \mathbf{U} \left( \mathbf{r}_v^t \odot \mathbf{h}_v^{t-1} \right) \right) \tag{16}$$

$$\mathbf{h}_v^t = \left( 1 - \mathbf{z}_v^t \right) \odot \mathbf{h}_v^{t-1} + \mathbf{z}_v^t \odot \widetilde{\mathbf{h}_v^t} \tag{17}$$

Where $\odot$ is element-wise multiplication, and $\mathbf{z}_v$ and $\mathbf{r}_v$ are the update and reset vectors, respectively
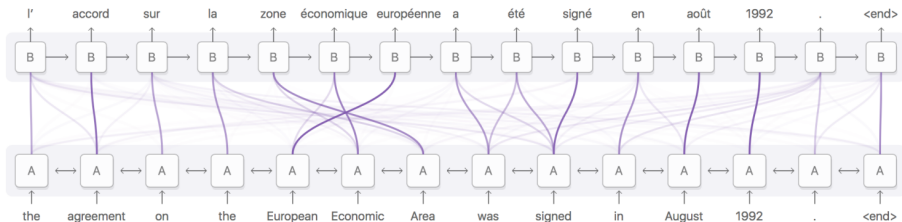
# Outline

- Gating mechanisms are designed for data that changes sequentially; however, our graphs have static features
- GAT incorporates attention mechanisms into the propagation step

# Attention

# Attention



GAT uses a *self-attention* mechanism to compute the hidden states of each node by attending over itself and its neighbors
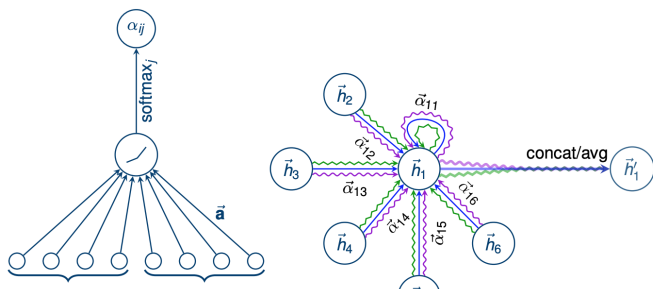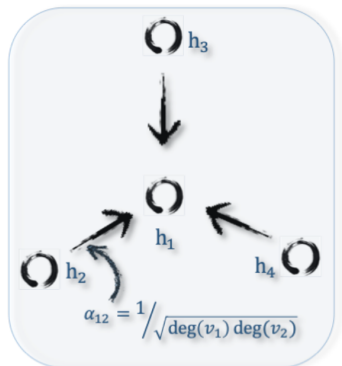
$$e_{vj} = a(\mathbf{W}^a \mathbf{h}_v, \mathbf{W}^b \mathbf{h}_j) \tag{18}$$

$$\alpha_{vj} = \frac{\exp(e_{vj})}{\sum_{k \in \mathcal{N}_v} \exp(e_{ik})} \tag{19}$$
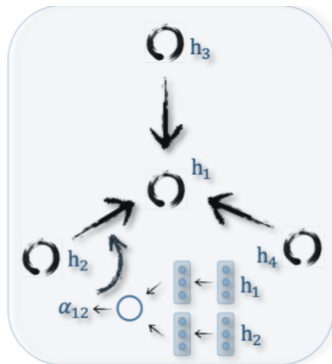
$$\mathbf{h}_v = \sigma\left( \sum_{j \in \mathcal{N}_v} \alpha_{vj} \mathbf{W}^z \mathbf{h}_j \right) \tag{20}$$

# Attention



(a) Graph Convolution Networks [14] explicitly assign a non-parametric weight $a_{ij} = \frac{1}{\sqrt{deg(v_i)deg(v_j)}}$ to the neighbor $v_j$ of $v_i$ during the aggregation process.

(b) Graph Attention Networks [15] implicitly capture the weight $a_{ij}$ via an end to end neural network architecture, so that more important nodes receive larger weights.

# Gating and Attention

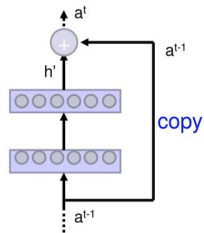| Name | Variant | Aggregator | Updater |
|---|---|---|---|
| Graph Attention Networks | GAT | $\alpha_{vk} = \dfrac{\exp\left(\text{LeakyReLU}\left(\mathbf{a}^T[\mathbf{W}\mathbf{h}_v \| \mathbf{W}\mathbf{h}_k]\right)\right)}{\sum_{j \in \mathcal{N}_v} \exp\left(\text{LeakyReLU}(\mathbf{a}^T[\mathbf{W}\mathbf{h}_v \| \mathbf{W}\mathbf{h}_j])\right)}$ <br> $\mathbf{h}_{\mathcal{N}_v}^t = \sigma\left(\sum_{k \in \mathcal{N}_v} \alpha_{vk}\mathbf{W}\mathbf{h}_k\right)$ <br> Multi-head concatenation: <br> $\mathbf{h}_{\mathcal{N}_v}^t = \Big\|_{m=1}^{M} \sigma\left(\sum_{k \in \mathcal{N}_v} \alpha_{vk}^m \mathbf{W}^m \mathbf{h}_k\right)$ <br> Multi-head average: <br> $\mathbf{h}_{\mathcal{N}_v}^t = \sigma\left(\frac{1}{M}\sum_{m=1}^{M}\sum_{k \in \mathcal{N}_v} \alpha_{vk}^m \mathbf{W}^m \mathbf{h}_k\right)$ | $\mathbf{h}_v^t = \mathbf{h}_{\mathcal{N}_v}^t$ |
| Gated Graph Neural Networks | GGNN | $\mathbf{h}_{\mathcal{N}_v}^t = \sum_{k \in \mathcal{N}_v} \mathbf{h}_k^{t-1} + \mathbf{b}$ | $\mathbf{z}_v^t = \sigma(\mathbf{W}^z \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{U}^z \mathbf{h}_v^{t-1})$ <br> $\mathbf{r}_v^t = \sigma(\mathbf{W}^r \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{U}^r \mathbf{h}_v^{t-1})$ <br> $\widetilde{\mathbf{h}_v^t} = \tanh(\mathbf{W}\mathbf{h}_{\mathcal{N}_v}^t + \mathbf{U}(\mathbf{r}_v^t \odot \mathbf{h}_v^{t-1}))$ <br> $\mathbf{h}_v^t = (1 - \mathbf{z}_v^t) \odot \mathbf{h}_v^{t-1} + \mathbf{z}_v^t \odot \widetilde{\mathbf{h}_v^t}$ |

# Outline

Uses uses layer-wise gates. The output of a layer is summed with its input with gating weights (inspired by Highway nets)

$$
\begin{aligned}
\mathbf{T}(\mathbf{h}^t) &= \sigma\left(\mathbf{W}^t\mathbf{h}^t + \mathbf{b}^t\right) \\
\mathbf{h}^{t+1} &= \mathbf{h}^{t+1} \odot \mathbf{T}(\mathbf{h}^t) + \mathbf{h}^t \odot (1 - \mathbf{T}(\mathbf{h}^t))
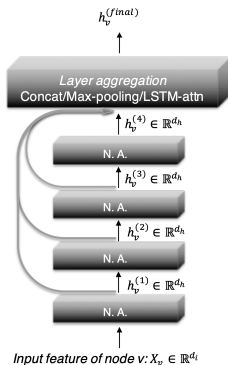\end{aligned}
\tag{21}
$$

- With neighborhood aggregation, the receptive field of each node grows exponentially w.r.t. the number of layers (steps) $T$
- The Jump Knowledge Network selects from all of the intermediate representations for each node at the last layer
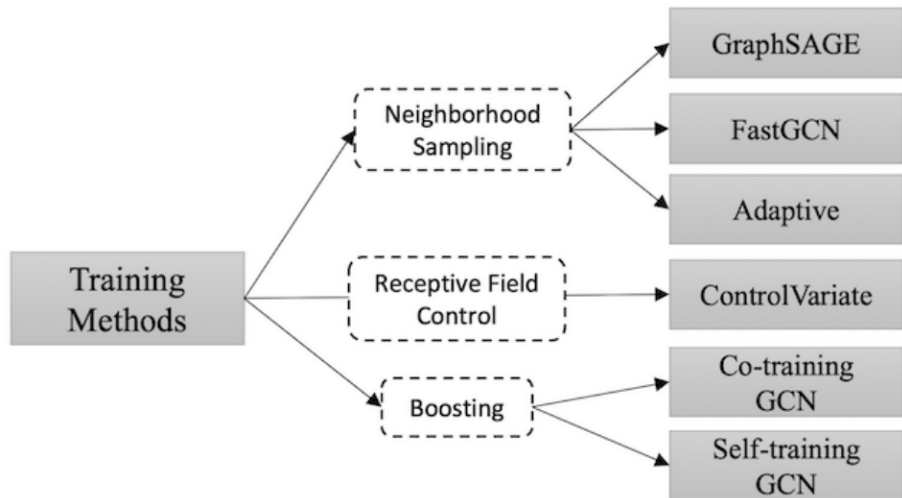  - Allows the model adapt the effective neighborhood size for each node as needed

$h_v^{(final)}$

Layer aggregation
Concat/Max-pooling/LSTM-attn

$h_v^{(4)} \in \mathbb{R}^{d_h}$

N. A.

$h_v^{(3)} \in \mathbb{R}^{d_h}$

N. A.

$h_v^{(2)} \in \mathbb{R}^{d_h}$

N. A.

$h_v^{(1)} \in \mathbb{R}^{d_h}$

N. A.

Input feature of node v: $X_v \in \mathbb{R}^{d_i}$

# Outline

- GraphSAGE solved the problems of the original GCN by replacing full graph Laplacian with learnable aggregation functions, which are key to generalize to unseen nodes.
- In addition, GraphSAGE uses neighbor sampling to alleviate receptive field expansion

- Chen et. al (2018) proposed a control-variate based stochastic approximation algorithms by utilizing the historical activations of nodes as a control variate.
- This limits the receptive field in the 1-hop neighborhood, but is efficient

- Li et. al. (2018) note that GNNs requires many additional labeled data for validation and also suffers from the localized nature of the convolutional filter
- To solve the limitations, the authors propose a method to find the nearest neighbors of training data and a boosting-like method
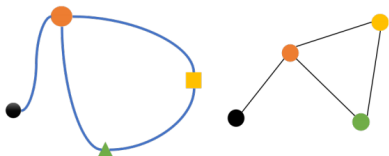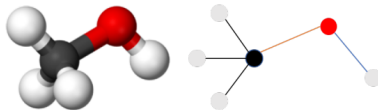
# Outline
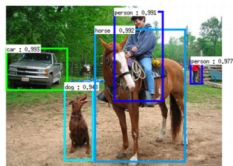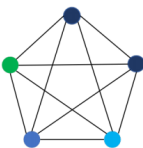
(a) physics



(b) molecule



(c) image



(d) text

# Commonly Used Datasets

| Category | Dataset | Source | # Graphs | # Nodes | # Edges | #Features | # Labels |
|---|---|---|---|---|---|---|---|
| Citation Networks | Cora | [103] | 1 | 2708 | 5429 | 1433 | 7 |
| | Citeseer | [103] | 1 | 3327 | 4732 | 3703 | 6 |
| | Pubmed | [103] | 1 | 19717 | 44338 | 500 | 3 |
| | DBLP | dblp.uni-trier.de [105](aminer.org /citation) | 1 | - | - | - | - |
| Social Networks | BlogCatalog | [107] | 1 | 10312 | 333983 | - | 39 |
| | Reddit | [24] | 1 | 232965 | 11606919 | 602 | 41 |
| | Epinions | www.epinions.com | 1 | - | - | - | - |
| Chemical/ Biological Graphs | PPI | [109] | 24 | 56944 | 818716 | 50 | 121 |
| | NCI-1 | [110] | 4100 | - | - | 37 | 2 |
| | NCI-109 | [110] | 4127 | - | - | 38 | 2 |
| | MUTAG | [111] | 188 | - | - | 7 | 2 |
| | D&D | [112] | 1178 | - | - | - | 2 |
| | QM9 | [113] | 133885 | - | - | - | 13 |
| | tox21 | tripod.nih.gov/ tox21/challenge/ | 12707 | - | - | - | 12 |
| Unstruct- ured Graphs | MNIST | yann.lecun.com /exdb/mnist/ | 70000 | - | - | - | 10 |
| | Wikipedia | www.mattmahoney .net/dc/textdata | 1 | 4777 | 184812 | - | 40 |
| | 20NEWS | [114] | 1 | 18846 | - | - | 20 |