# How Powerful Are Graph Neural Networks

Keyulu Xu, Weihua Hu, Jure Leskovec, Stefanie Jegelka

Presented by Mikhail Mishin[1]

[1]Aalto University
mikhail.mishin@aalto.fi

2019

# Agenda

- Intro to Graph Neural Networks (GNNs)
- Theoretical framework: Characterize GNN's discriminative power
- Propose a Maximally Powerful GNN
- Experiments
- Summary & Conclusion
- Demo & Discussion

# Setup

Assume we have a graph $G = (V, E)$

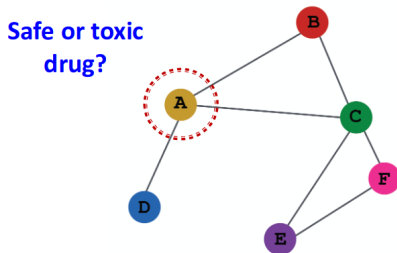- $N(v)$ is the neighborhood of node $v$ (the set of nodes adjacent to $v$)
- $X_v$ is the input feature vector of node $v$
  - Social networks: user profile, user image
  - Biological networks: gene expression profile
  - No features: node degree, constant
- Assume the node input features are from a countable universe
- Can assign each feature vector a unique label in $\{a, b, c, \dots\}$
- Feature vectors of a set of neighbouring nodes form a multiset (allows multiple instances)

# Supervised Learning on Graphs

**Node classification**

▶ Given a label $y_v$ for each node $v \in V$

▶ Learn an embedding $h_v$ of $v$

▶ To help predict $v$'s label, $y_v = f(h_v)$
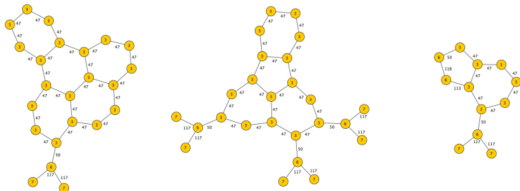
Example: a drug interaction network


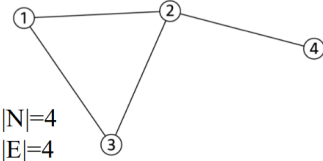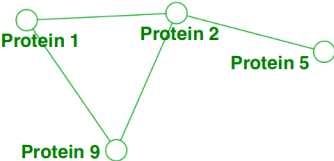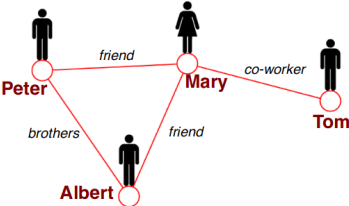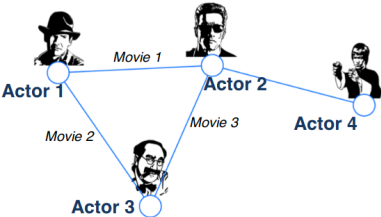
**Safe or toxic drug?**

# Supervised Learning on Graphs

## Graph classification

- Given a set of graphs $\{G_1, \ldots, G_N\} \subseteq \mathcal{G}$
- And their labels $\{y_1, \ldots, y_N\} \subseteq \mathcal{Y}$
- Learn an embedding $h_G$ of a graph $G$
- To help predict $G$'s label, $y_G = g(h_G)$

Example: aromatic and heteroaromatic nitro compounds, MUTAG
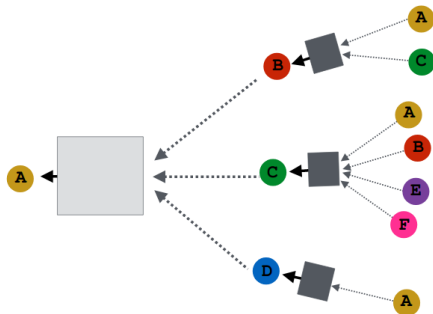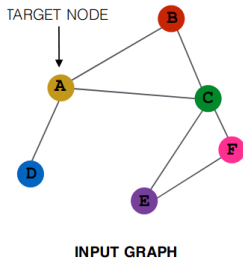
# More Graph Examples

# Graph Neural Networks[1]

**Key idea:** generate node embeddings based on local network neighborhoods



---

[1]Leskovec, *Stanford CS224W: Machine Learning with Graphs.*

# Graph Neural Networks[2]

Intuition: nodes aggregate information from their neighbors using neural networks



**Neural networks**

[2]Leskovec, *Stanford CS224W: Machine Learning with Graphs.*

# Graph Neural Networks[3]

**Intuition:** network neighborhood defines a computational graph

Every node defines a computation graph based on its neighborhood!



INPUT GRAPH

# Graph Neural Networks[4]

- ▶ Model can be of arbitrary depth
- ▶ Nodes have embeddings at each layer
- ▶ Layer-0 embedding of node $v$ is its input feature vector $X_v$
- ▶ Layer-$k$ embedding gets information that's $k$ hops away



[4]Leskovec, *Stanford CS224W: Machine Learning with Graphs.*

# Neighborhood Aggregation[5]

- ▶ Key distinctions are in how different approaches aggregate information across layers
- ▶ The same aggregation parameters are shared for each layer



What is in the box?

TARGET NODE

INPUT GRAPH

[5]Leskovec, *Stanford CS224W: Machine Learning with Graphs*.

# Neighborhood Aggregation

- Formally, the $k$-th layer of a GNN is:

$$a_v^{(k)} = \text{AGGREGATE}^{(k)} \left( \left\{ h_u^{(k-1)} : u \in N(v) \right\} \right)$$

$$h_v^{(k)} = \text{COMBINE}^{(k)} \left( h_v^{(k-1)}, a_v^{(k)} \right)$$

- $h_v^{(k)}$ is the embedding of node $v$ at the $k$-th iteration/layer.
- $h_v^{(0)} = X_v$, $v$'s input feature vector

# GNN Architectures: GraphSAGE[6]

▶ AGGREGATE is formulated as:

$$a_v^{(k)} = \text{MAX}\left(\left\{\text{ReLU}\left(W_{\text{pool}}^{(k)} \cdot h_u^{(k-1)}\right), \forall u \in N(v)\right\}\right)$$

▶ COMBINE is formulated as:

$$h_v^{(k)} = W^{(k)} \cdot \text{CONCAT}\left[h_v^{(k-1)}, a_v^{(k)}\right]$$

▶ MAX represents the element-wise max pooling operation

▶ $W_{\text{pool}}^{(k)}$, $W^{(k)}$ are trainable matrices

---

[6]Hamilton, Ying, and Leskovec, "Inductive Representation Learning on Large Graphs".

- AGGREGATE and COMBINE are formulated as:

$$h_v^{(k)} = \text{ReLU}\left(W^{(k)} \cdot \text{MEAN}\left\{h_u^{(k-1)}, \forall u \in N(v) \cup \{v\}\right\}\right)$$

- MEAN represents the element-wise mean pooling operation
- $W^{(k)}$ is a trainable matrix

---

[7] Kipf and Welling, "Semi-Supervised Classification with Graph Convolutional Networks".

# Supervised Training[8]

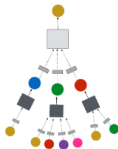Node classification: use the embedding $h_v^{(K)}$ of the final iteration for prediction

$$\mathcal{L} = \sum_{v \in V} y_v \log(\sigma(\mathbf{z}_v^\top \boldsymbol{\theta})) + (1 - y_v) \log(1 - \sigma(\mathbf{z}_v^\top \boldsymbol{\theta}))$$

**Encoder output:** node embedding

Classification weights

Node class label

Safe or toxic drug?



[8]Leskovec, *Stanford CS224W: Machine Learning with Graphs.*

# Supervised Training

Graph classification:

- Aggregate node embeddings from the final iteration with a READOUT function:
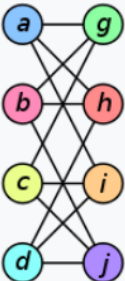
$$h_G = \text{READOUT}\left( \left\{ h_v^{(K)} | v \in V \right\} \right)$$

- E.g. READOUT = SUM or READOUT = MEAN
- Can be a more sophisticated pooling function
- To train use $h_G$ for prediction, same as for the node classification task

# Graph Isomorphism Problem

- ▶ Are two graphs topologically identical?
- ▶ Known to be in NP, no polynomial-time algorithm
- ▶ GNN is able to map two different graphs to different embeddings $\implies$ it solves GI
- ▶ Might need a weaker criterion...

| Graph G | Graph H | An isomorphism between G and H |
|---|---|---|
|  |  | $f(a) = 1$ <br> $f(b) = 6$ <br> $f(c) = 8$ <br> $f(d) = 3$ <br> $f(g) = 5$ <br> $f(h) = 2$ <br> $f(i) = 4$ <br> $f(j) = 7$ |

# Weisfeiler-Lehman Test of Graph Isomorphism

- The 1-dimensional WL test iteratively:
  - Aggregates the labels of nodes and their neighborhoods
  - Hashes the aggregated labels into unique new labels



- Decide two graphs are non-isomorphic if at some iteration the multisets of labels between the two graphs differ

# GNNs are At Most As Powerful As WL

- **Key contribution 1:** GNNs are at most as powerful as the WL test in distinguishing graph structures.
- GNN maps $G_1$ and $G_2$ to different embeddings $\implies$ the WL test decides $G_1$ and $G_2$ are non-isomorphic
- See the formal proof is in the paper (Lemma 2)



2 WL test iterations    Captures structures

**Multiset**

**Graph**                    Rooted subtree                    GNN aggregation

# Maximally Powerful GNNs (MP-GNNs)

▶ A Maximally Powerful GNN would never map two different node neighborhoods to the same representation

▶ Its AGGREGATE, COMBINE, READOUT must be injective



**Graph**      Rooted subtree      GNN aggregation

# MP-GNNs are As Powerful As WL

- **Key contribution 2:** MP-GNNs are as powerful as the WL test in distinguishing graph structures.
- The WL test decides $G_1$ and $G_2$ are non-isomorphic $\implies$ a MP-GNN maps $G_1$ and $G_2$ to different embeddings
- Note, the proof only holds for countable sets of input node features
- See the formal proof in the paper (Theorem 3 and Lemma 4)



**Graph**        2 WL test iterations        Rooted subtree        Captures structures        **Multiset**        GNN aggregation

# GNNs Learn to Embed

- WL test can only discriminate different graph structures
- GNNs learn useful node representations, capturing similarity of graph structures
- WL test provides theoretical context for GNN design

# GraphSAGE is Not Maximally Powerful

▶ AGGREGATE is formulated as:

$$a_v^{(k)} = \text{MAX}\left(\left\{\text{ReLU}\left(W_{\text{pool}}^{(k)} \cdot h_u^{(k-1)}\right), \forall u \in N(v)\right\}\right)$$

▶ COMBINE is formulated as:

$$h_v^{(k)} = W^{(k)} \cdot \text{CONCAT}\left[h_v^{(k-1)}, a_v^{(k)}\right]$$

▶ MAX is not injective

# Max-Pooling Learns Sets with Distinct Elements

▶ MAX aggregator ignores multiplicities



**Input**   max - set

▶ MAX fails



vs.   vs.

# GCN is Not Maximally Powerful

- ▶ AGGREGATE and COMBINE are formulated as:

$$h_v^{(k)} = \text{ReLU}\left(W^{(k)} \cdot \text{MEAN}\left\{h_u^{(k-1)}, \forall u \in N(v) \cup \{v\}\right\}\right)$$

- ▶ MEAN is not injective

# Mean-Pooling Learns Distributions

▶ MEAN captures the proportion of elements of a given type



▶ MEAN fails

# Sum is Injective

► Key contribution 3: Ranking by representational power

► SUM captures the whole multiset



Input — sum - multiset > mean - distribution > max - set

► SUM succeeds



(a) Mean and Max both fail    (b) Max fails    (c) Mean and Max both fail

# Graph Isomorphism Network (GIN)

- **Key contribution 4:** GIN is provably maximally powerful
- **AGGREGATE** and **COMBINE** are formulated as:

$$h_v^{(k)} = \text{MLP}^{(k)} \left( \left( 1 + \epsilon^{(k)} \right) \cdot h_v^{(k-1)} + \sum_{u \in N(v)} h_u^{(k-1)} \right)$$

- MLP is Multi-Layer Perceptron
- $\epsilon$ is needed to distinguish the root/central node
- 1-layer perceptron is not enough
- See Lemma 5 and Corollary 6

# Graph-Level READOUT of GIN

- Concatenate graph representations across all layers of GIN

$$h_G = \text{CONCAT}\left(\text{READOUT}\left(\left\{h_v^{(k)}|v \in V\right\}\right)|k = 0, 1, \ldots K\right)$$

- Node embeddings get more refined and global as the number layers increases
- Earlier layers may generalize better
- If READOUT=SUM, provably maximally powerful
- However READOUT=MEAN, performs better on some datasets

# Experiments: Datasets

- **The goal** is to learn from the network structure, not to rely on input features
- **Bioinformatic datasets:** MUTAG, PRC, NCI1, PROTEINS
  - Graphs represent chemical compounds
  - MUTAG: aromatic and heteroaromatic nitro compounds
  - Nodes have categorial input features
- **Social network datasets:** COLLAB, IMDB-BINARY, IMDB-MULTI, REDDIT-BINARY, REDDIT-MULTI5K
  - Graphs represent social communities
  - IMDB: ego-network for each actor/actress, classify genre
  - REDDIT: online discussion, classify subreddit
  - Node input features: IMDB - node degrees, REDDIT - constant

# Experiments: Models and Configurations

- Evaluate GIN:
  - GIN-$\epsilon$ with learning $\epsilon$
  - GIN-0 without learning $\epsilon$
- Evaluate less powerful versions:
  - Replace SUM with MEAN or MAX
  - Replace MLP with 1-layer perceptrons
- Apply the same graph-level READOUT:
  - SUM on bioinformatics datasets
  - MEAN on social datasets (better perfomance on test)
- For all configurations, apply 5 GNN layers (including input)
- All MLPs have 2 layers
- Apply Batch normalization on hidden layers
- Compare with state-of-the-art baselines: PATCHY-SAN, Deep Graph CNN, etc.

# Experiments: Training Performance



- ▶ GIN-$\epsilon$ and GIN-0 almost perfectly fit training
- ▶ Others slightly underfit
- ▶ Observed training accuracy aligns with representational power
- ▶ WL subtree kernel has the best training accuracies

# Experiments: Test Performance

| | Datasets | IMDB-B | IMDB-M | RDT-B | RDT-M5K | COLLAB | MUTAG | PROTEINS | PTC | NCI1 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Datasets** | # graphs | 1000 | 1500 | 2000 | 5000 | 5000 | 188 | 1113 | 344 | 4110 |
| | # classes | 2 | 3 | 2 | 5 | 3 | 2 | 2 | 2 | 2 |
| | Avg # nodes | 19.8 | 13.0 | 429.6 | 508.5 | 74.5 | 17.9 | 39.1 | 25.5 | 29.8 |
| **Baselines** | WL subtree | 73.8 ± 3.9 | 50.9 ± 3.8 | 81.0 ± 3.1 | 52.5 ± 2.1 | 78.9 ± 1.9 | 90.4 ± 5.7 | 75.0 ± 3.1 | 59.9 ± 4.3 | **86.0 ± 1.8** * |
| | DCNN | 49.1 | 33.5 | – | – | 52.1 | 67.0 | 61.3 | 56.6 | 62.6 |
| | PATCHYSAN | 71.0 ± 2.2 | 45.2 ± 2.8 | 86.3 ± 1.6 | 49.1 ± 0.7 | 72.6 ± 2.2 | **92.6 ± 4.2** * | 75.9 ± 2.8 | 60.0 ± 4.8 | 78.6 ± 1.9 |
| | DGCNN | 70.0 | 47.8 | – | – | 73.7 | 85.8 | 75.5 | 58.6 | 74.4 |
| | AWL | 74.5 ± 5.9 | 51.5 ± 3.6 | 87.9 ± 2.5 | 54.7 ± 2.9 | 73.9 ± 1.9 | 87.9 ± 9.8 | – | – | – |
| **GNN variants** | SUM–MLP (**GIN-0**) | **75.1 ± 5.1** | **52.3 ± 2.8** | **92.4 ± 2.5** | **57.5 ± 1.5** | **80.2 ± 1.9** | **89.4 ± 5.6** | **76.2 ± 2.8** | **64.6 ± 7.0** | **82.7 ± 1.7** |
| | SUM–MLP (**GIN-$\epsilon$**) | **74.3 ± 5.1** | **52.1 ± 3.6** | **92.2 ± 2.3** | **57.0 ± 1.7** | **80.1 ± 1.9** | **89.0 ± 6.0** | **75.9 ± 3.8** | 63.7 ± 8.2 | **82.7 ± 1.6** |
| | SUM–1-LAYER | 74.1 ± 5.0 | **52.2 ± 2.4** | 90.0 ± 2.7 | 55.1 ± 1.6 | **80.6 ± 1.9** | **90.0 ± 8.8** | **76.2 ± 2.6** | 63.1 ± 5.7 | 82.0 ± 1.5 |
| | MEAN–MLP | 73.7 ± 3.7 | **52.3 ± 3.1** | 50.0 ± 0.0 | 20.0 ± 0.0 | 79.2 ± 2.3 | 83.5 ± 6.3 | 75.5 ± 3.4 | **66.6 ± 6.9** | 80.9 ± 1.8 |
| | MEAN–1-LAYER (GCN) | 74.0 ± 3.4 | 51.9 ± 3.8 | 50.0 ± 0.0 | 20.0 ± 0.0 | 79.0 ± 1.8 | 85.6 ± 5.8 | 76.0 ± 3.2 | 64.2 ± 4.3 | 80.2 ± 2.0 |
| | MAX–MLP | 73.2 ± 5.8 | 51.1 ± 3.6 | – | – | – | 84.0 ± 6.1 | 76.0 ± 3.2 | 64.6 ± 10.2 | 77.8 ± 1.3 |
| | MAX–1-LAYER (GraphSAGE) | 72.3 ± 5.3 | 50.9 ± 2.2 | – | – | – | 85.1 ± 7.6 | 75.9 ± 3.2 | 63.9 ± 7.7 | 77.7 ± 1.5 |

▶ GINs outperform less powerful GNN varians

▶ GIN-0 slightly outperforms GIN-$\epsilon$

▶ GINs shine on social network datasets

▶ MEAN GNNs fail to capture any structures of the unlabeled datasets (RDT-B, RDT-M5K)

# TL;DR:

- POWER - ability to discriminate graph structures
- Prove $POWER(GNNs) \leq POWER(WL)$
- Establish conditions for $POWER(GNNs) = POWER(WL)$
- Show $POWER(GraphSAGE) < POWER(WL)$ and $POWER(GCN) < POWER(WL)$
- Propose GIN and prove $POWER(GIN) = POWER(WL)$
- Empirically show GIN performs better than GraphSAGE and GCN

# Conclusion[9]

Graph Convolutional Neural Networks:

- ► Representation learning paradigm can be extended to graphs
- ► Can effectively combine node attribute data with the network information
- ► State-of-the-art results in a number of domains/tasks
- ► Use end-to-end training instead of multi-stage approaches for better performance

---

[9]Leskovec, *Stanford CS224W: Machine Learning with Graphs.*