

# EECE 571F: Deep Learning with Structures

## Lecture 6: Unsupervised/Self-supervised Graph Representation Learning

Renjie Liao

University of British Columbia

Winter, Term 1, 2022

# Course Scope

- Brief Intro to Deep Learning
- Geometric Deep Learning
  - Deep Learning Models for Sets and Sequences: Deep Sets & Transformers
  - Deep Learning Models for Graphs: Graph Convolution & Message Passing GNNs
  - Expressiveness & Generalizations of GNNs
  - Unsupervised/Self-supervised Graph Representation Learning
- Probabilistic Deep Learning
  - Deep Generative Models:  
Auto-regressive models, GANs, VAEs, Diffusion/Score based models
  - Discrete/Hybrid Latent Variable Models: RBMs, Latent Graph Models
  - Stochastic Gradient Estimation

# Course Scope

- Brief Intro to Deep Learning
- Geometric Deep Learning
  - Deep Learning Models for Sets and Sequences: Deep Sets & Transformers
  - Deep Learning Models for Graphs: Graph Convolution & Message Passing GNNs
  - Expressiveness & Generalizations of GNNs
  - **Unsupervised/Self-supervised Graph Representation Learning**
- Probabilistic Deep Learning
  - Deep Generative Models:  
Auto-regressive models, GANs, VAEs, Diffusion/Score based models
  - Discrete/Hybrid Latent Variable Models: RBMs, Latent Graph Models
  - Stochastic Gradient Estimation

# Unsupervised Graph Representation Learning

- Deep Generative Models of Graphs

Building probabilistic distributions of graphs (e.g., adjacency matrix  $A$ ) and node representations  $X$

# Unsupervised Graph Representation Learning

- Deep Generative Models of Graphs

Building probabilistic distributions of graphs (e.g., adjacency matrix  $A$ ) and node representations  $X$

- In this lecture, we focus on methods that learn good node representations

Given graphs (e.g., adjacency matrix  $A$ ), the goal is to learn node representations  $X$

The learned  $X$  is useful for supervised fine-tuning, e.g., node classification

# Unsupervised / Self-Supervised Learning

Since only data is given, we need a learning criterion:

# Unsupervised / Self-Supervised Learning on Graphs

Since only data is given, we need a learning criterion:

- Likelihood (Autoregressive models)
- Reconstruction Loss (Auto-encoders)
- Contrastive Loss (Noise contrastive estimation, Self-supervised learning)
- Min-max Loss (Generative adversarial networks)

.....

# Unsupervised / Self-Supervised Learning on Graphs

Since only data is given, we need a learning criterion:

- **Likelihood (Autoregressive models)**
- Reconstruction Loss (Auto-encoders)
- Contrastive Loss (Noise contrastive estimation, Self-supervised learning)
- Min-max Loss (Generative adversarial networks)

.....

DeepWalk [1] Node2Vec[2]



# Unsupervised / Self-Supervised Learning on Graphs

Since only data is given, we need a learning criterion:

- **Likelihood (Autoregressive models)** DeepWalk [1] Node2Vec[2]
- Reconstruction Loss (Auto-encoders)
- **Contrastive Loss (Noise contrastive estimation, Self-supervised learning)** DeepGraphInfoMax [3]
- **Min-max Loss (Generative adversarial networks)** DeepGraphInfoMax [3]

.....

# Language Models

Model the probability of words given its context (e.g., a fixed-size window):

$$p(w_{i+1} | w_i, \dots, w_{i-K})$$

# Language Models

Model the probability of words given its context (e.g., a fixed-size window):

$$p(w_{i+1} | w_i, \dots, w_{i-K})$$

$$p(w_i | \{w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k}\})$$

# Language Models

Model the probability of words given its context (e.g., a fixed-size window):

$$p(w_{i+1} | w_i, \dots, w_{i-K})$$

$$p(w_i | \{w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k}\})$$

- One can use RNNs or CNNs to construct the probability
- The model can be learned via *maximum likelihood*

# Skip-Gram Models

Model the probability of its context (e.g., a fixed-size window) given the word:

$$p(\{w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k}\} | w_i) = \prod_{j \neq i, j \in \mathcal{N}_i} p(w_j | w_i)$$

# Skip-Gram Models

Model the probability of its context (e.g., a fixed-size window) given the word:

$$p(\{w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k}\} | w_i) = \prod_{j \neq i, j \in \mathcal{N}_i} p(w_j | w_i)$$

- Assumes conditional independence, i.e., the context is orderless
- The model only takes one word as input, thus being efficient
- Interpolation (context  $\Rightarrow$  word) vs. extrapolation (word  $\Rightarrow$  context)

# Skip-Gram Models

Model the probability of its context (e.g., a fixed-size window) given the word:

$$p(\{w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k}\} | w_i) = \prod_{j \neq i, j \in \mathcal{N}_i} p(w_j | w_i)$$

- Assumes conditional independence, i.e., the context is orderless
- The model only takes one word as input, thus being efficient
- Interpolation (context  $\Rightarrow$  word) vs. extrapolation (word  $\Rightarrow$  context)

Can we generalize this model to graphs?

# Skip-Gram Models

Model the probability of its context (e.g., a fixed-size window) given the word:

$$p(\{w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k}\} | w_i) = \prod_{j \neq i, j \in \mathcal{N}_i} p(w_j | w_i)$$

- Assumes conditional independence, i.e., the context is orderless
- The model only takes one word as input, thus being efficient
- Interpolation (context  $\Rightarrow$  word) vs. extrapolation (word  $\Rightarrow$  context)

Words  $\rightarrow$  Nodes

Can we generalize this model to graphs?



# Skip-Gram Models

Model the probability of its context (e.g., a fixed-size window) given the word:

$$p(\{w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k}\} | w_i) = \prod_{j \neq i, j \in \mathcal{N}_i} p(w_j | w_i)$$

- Assumes conditional independence, i.e., the context is orderless
- The model only takes one word as input, thus being efficient
- Interpolation (context  $\Rightarrow$  word) vs. extrapolation (word  $\Rightarrow$  context)

Words  $\rightarrow$  Nodes

Can we generalize this model to graphs?

How about edges?

# Random Walks

A special Markov Chain

*Starting at a node, one can randomly choose a neighboring node at a time to walk*

# Random Walks

A special Markov Chain

*Starting at a node, one can randomly choose a neighboring node at a time to walk*

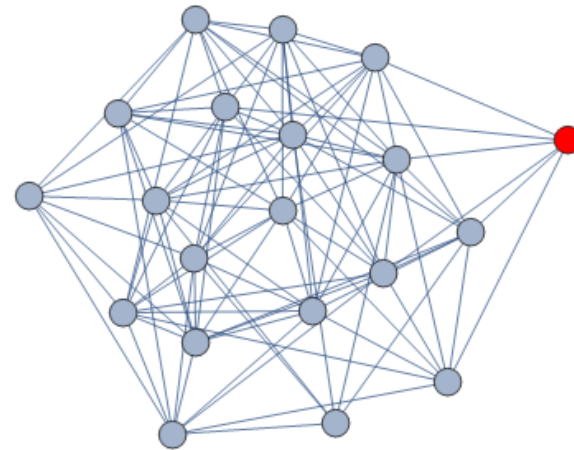
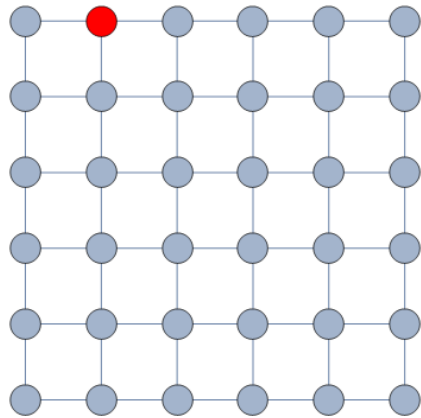
$$p_{ij} = \frac{1}{|\mathcal{N}_i|}$$

# Random Walks

A special Markov Chain

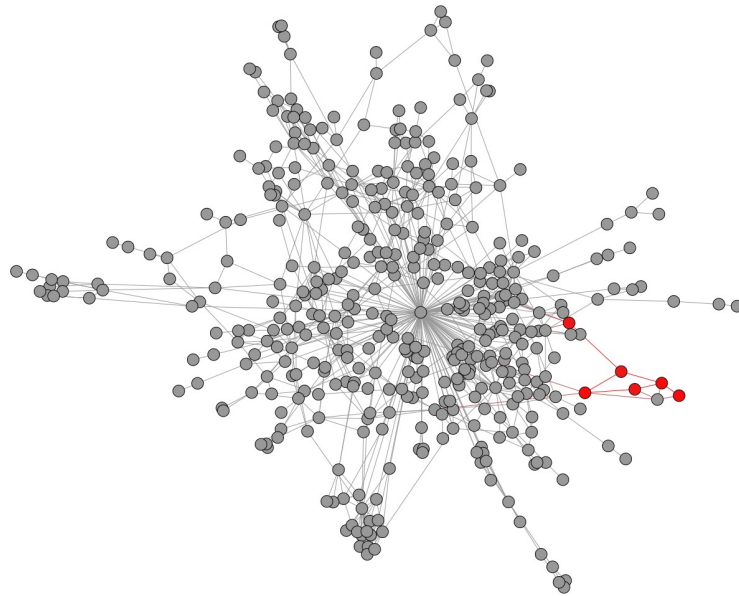
*Starting at a node, one can randomly choose a neighboring node at a time to walk*

$$p_{ij} = \frac{1}{|\mathcal{N}_i|}$$

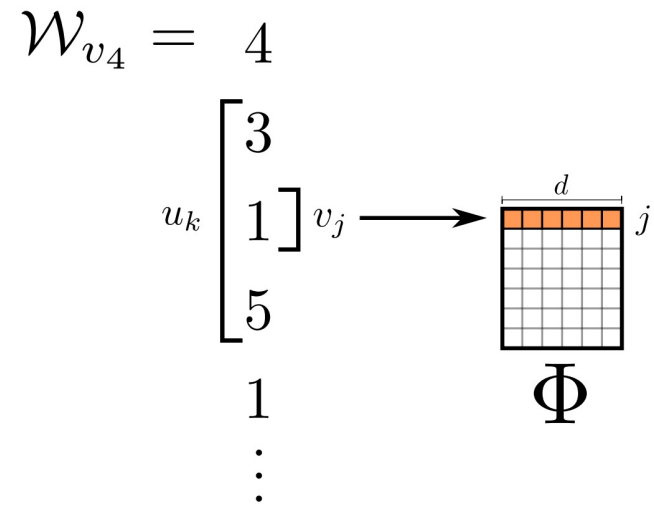


# Deep Walks

Model the probability of context (random walks) given its word (vertex):



(a) Random walk generation.

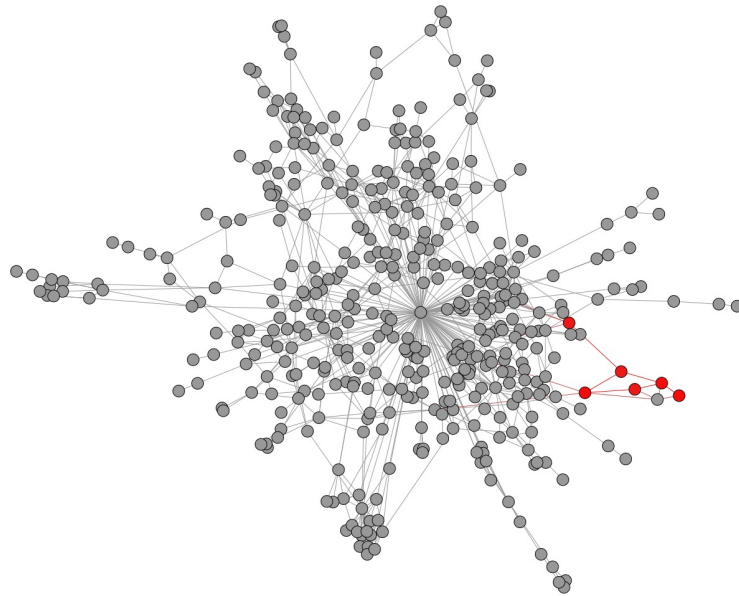


(b) Representation mapping.

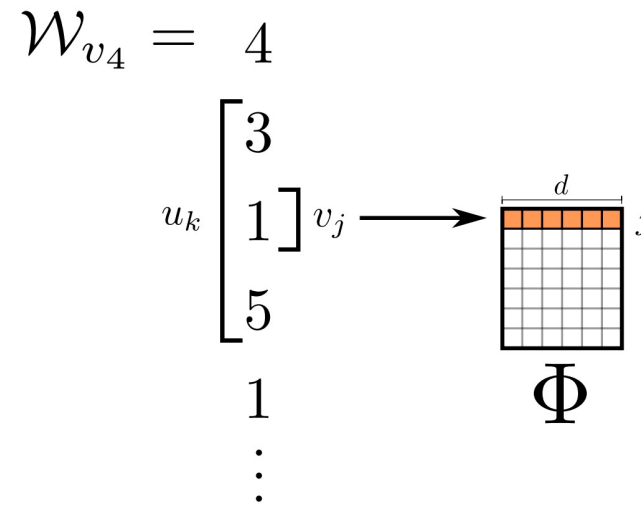
# Deep Walks

Model the probability of context (random walks) given its word (vertex):

$$\max \log p(\{v_{i-w}, \dots, v_{i-1}, v_{i+1}, \dots, v_{i+w}\} \mid \Phi(v_i))$$



(a) Random walk generation.



(b) Representation mapping.

# Deep Walks

Model the probability of context (random walks) given its word (vertex):

---

**Algorithm 1** DEEPWALK( $G, w, d, \gamma, t$ )

---

**Input:** graph  $G(V, E)$

    window size  $w$

    embedding size  $d$

    walks per vertex  $\gamma$

    walk length  $t$

**Output:** matrix of vertex representations  $\Phi \in \mathbb{R}^{|V| \times d}$

1: Initialization: Sample  $\Phi$  from  $\mathcal{U}^{|V| \times d}$

2: Build a binary Tree  $T$  from  $V$

3: **for**  $i = 0$  to  $\gamma$  **do**

4:      $\mathcal{O} = \text{Shuffle}(V)$

5:     **for each**  $v_i \in \mathcal{O}$  **do**

6:          $\mathcal{W}_{v_i} = \text{RandomWalk}(G, v_i, t)$

7:          $\text{SkipGram}(\Phi, \mathcal{W}_{v_i}, w)$

8:     **end for**

9: **end for**

---

# Deep Walks

## Skip-Gram Algorithm

---

**Algorithm 2** SkipGram( $\Phi, \mathcal{W}_{v_i}, w$ )

---

```
1: for each  $v_j \in \mathcal{W}_{v_i}$  do  
2:   for each  $u_k \in \mathcal{W}_{v_i}[j - w : j + w]$  do  
3:      $J(\Phi) = -\log \Pr(u_k \mid \Phi(v_j))$   
4:      $\Phi = \Phi - \alpha * \frac{\partial J}{\partial \Phi}$   
5:   end for  
6: end for
```

---



# Deep Walks

## Skip-Gram Algorithm

---

**Algorithm 2** SkipGram( $\Phi, \mathcal{W}_{v_i}, w$ )

---

```
1: for each  $v_j \in \mathcal{W}_{v_i}$  do
2:   for each  $u_k \in \mathcal{W}_{v_i} [j - w : j + w]$  do
3:      $J(\Phi) = -\log \Pr(u_k | \Phi(v_j))$ 
4:      $\Phi = \Phi - \alpha * \frac{\partial J}{\partial \Phi}$ 
5:   end for
6: end for
```

---

## Conditional Independence

$$p(\{u_{j-w}, \dots, u_{j-1}, u_{j+1}, \dots, u_{j+w}\} | \Phi(v_j)) = \prod_{k \neq j, k \in \mathcal{N}_j} p(u_k | \Phi(v_j))$$

# Deep Walks

Model the probability of context (random walks) given its word (vertex)

Construct probability using Softmax:

$$p(u_k = i | \Phi(v_j)) = \frac{\exp(w_i^T \Phi(v_j))}{\sum_{m=1}^{|V|} \exp(w_m^T \Phi(v_j))}$$

# Deep Walks

Model the probability of context (random walks) given its word (vertex)

Construct probability using Softmax:

$$p(u_k = i | \Phi(v_j)) = \frac{\exp(w_i^T \Phi(v_j))}{\sum_{m=1}^{|V|} \exp(w_m^T \Phi(v_j))}$$

One weight per vertex, i.e., a huge softmax on large graphs

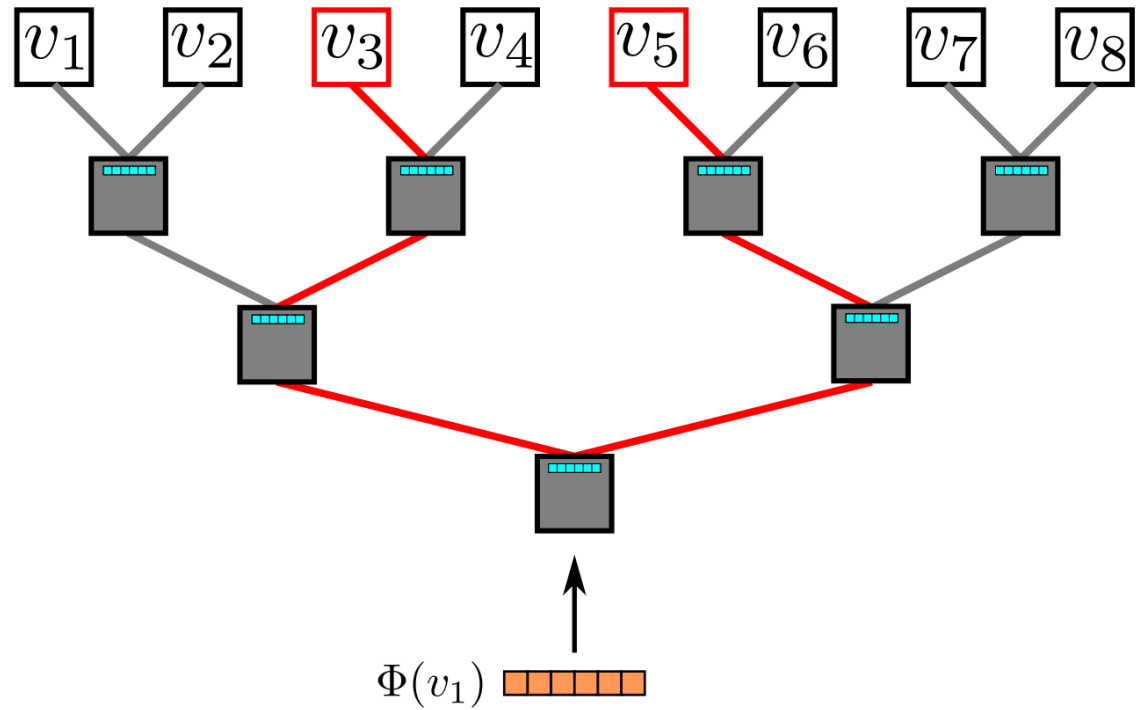
Can we improve the efficiency?

# Deep Walks

Hierarchical Softmax

Build a binary tree over vocabulary (the set of all vertices)

$$p(u_k | \Phi(v_j)) = \prod_{l=1}^{\lceil \log |V| \rceil} p(b_l | \Phi(v_j))$$



# Deep Walks

Multi-Label Classification (BlogCatalog)

Two-stage pipeline: 1) learn node embeddings unsupervisedly; 2) learn node classifier supervisedly

	% Labeled Nodes	10%	20%	30%	40%	50%	60%	70%	80%	90%
Micro-F1(%)	DEEPWALK	<b>36.00</b>	<b>38.20</b>	<b>39.60</b>	<b>40.30</b>	<b>41.00</b>	<b>41.30</b>	41.50	41.50	42.00
	SpectralClustering	31.06	34.95	37.27	38.93	39.97	40.99	<b>41.66</b>	<b>42.42</b>	<b>42.62</b>
	EdgeCluster	27.94	30.76	31.85	32.99	34.12	35.00	34.63	35.99	36.29
	Modularity	27.35	30.74	31.77	32.97	34.09	36.13	36.08	37.23	38.18
	wvRN	19.51	24.34	25.62	28.82	30.37	31.81	32.19	33.33	34.28
	Majority	16.51	16.66	16.61	16.70	16.91	16.99	16.92	16.49	17.26
Macro-F1(%)	DEEPWALK	<b>21.30</b>	<b>23.80</b>	25.30	26.30	27.30	27.60	27.90	28.20	28.90
	SpectralClustering	19.14	23.57	<b>25.97</b>	<b>27.46</b>	<b>28.31</b>	<b>29.46</b>	<b>30.13</b>	<b>31.38</b>	<b>31.78</b>
	EdgeCluster	16.16	19.16	20.48	22.00	23.00	23.64	23.82	24.61	24.92
	Modularity	17.36	20.00	20.80	21.85	22.65	23.41	23.89	24.20	24.97
	wvRN	6.25	10.13	11.64	14.24	15.86	17.18	17.98	18.86	19.57
	Majority	2.52	2.55	2.52	2.58	2.58	2.63	2.61	2.48	2.62

# Deep Graph InfoMax

Mutual Information between X and Y:

$$\begin{aligned} \text{MI}(X, Y) &= \text{KL}(p(X, Y) \| p(X)p(Y)) \\ &= \int \int p(X, Y) \log \frac{p(X, Y)}{p(X)p(Y)} dX dY \end{aligned}$$

# Deep Graph InfoMax

Mutual Information between X and Y:

$$\begin{aligned} \text{MI}(X, Y) &= \text{KL}(p(X, Y) \| p(X)p(Y)) \\ &= \int \int p(X, Y) \log \frac{p(X, Y)}{p(X)p(Y)} dX dY \end{aligned}$$

It quantifies the "amount of information" obtained about one random variable by observing the other random variable

The higher the mutual information => knowing one would give you more information about the other

# Deep Graph InfoMax

How to estimate Mutual Information?

$$\begin{aligned} \text{MI}(X, Y) &= \text{KL}(p(X, Y) \| p(X)p(Y)) \\ &= \int \int p(X, Y) \log \frac{p(X, Y)}{p(X)p(Y)} dX dY \\ &= \int \int p(X, Y) \log \frac{p(Y|X)}{p(Y)} dX dY \\ &= \mathbb{E}_{p(X, Y)} \left[ \log \frac{p(Y|X)}{p(Y)} \right] \end{aligned}$$



# Deep Graph InfoMax

How to estimate Mutual Information?

Suppose we generate  $Y$  from a mixture distribution:

$$q(C = 1) = q(C = 0) = \frac{1}{2}$$

$$p(Y|X) = q(Y|C = 1)$$

$$p(Y) = q(Y|C = 0)$$

# Deep Graph InfoMax

How to estimate Mutual Information?

Suppose we generate  $Y$  from a mixture distribution:

$$q(C = 1) = q(C = 0) = \frac{1}{2}$$

$$p(Y|X) = q(Y|C = 1)$$

$$p(Y) = q(Y|C = 0)$$

$$q(Y) = q(Y|C = 0)q(C = 0) + q(Y|C = 1)q(C = 1)$$

# Deep Graph InfoMax

How to estimate Mutual Information?

Suppose we generate  $Y$  from a mixture distribution:

$$q(C = 1) = q(C = 0) = \frac{1}{2}$$

$$p(Y|X) = q(Y|C = 1)$$

$$p(Y) = q(Y|C = 0)$$

$$q(Y) = q(Y|C = 0)q(C = 0) + q(Y|C = 1)q(C = 1)$$

Density Ratio Trick:

$$\frac{p(Y|X)}{p(Y)} = \frac{q(Y|C = 1)}{q(Y|C = 0)} = \frac{\frac{q(C=1|Y)q(Y)}{q(C=1)}}{\frac{q(C=0|Y)q(Y)}{q(C=0)}} = \frac{q(C = 1|Y)}{q(C = 0|Y)}$$

# Deep Graph InfoMax

How to estimate Mutual Information?

Suppose we generate  $Y$  from a mixture distribution:

$$q(C = 1) = q(C = 0) = \frac{1}{2}$$

$$p(Y|X) = q(Y|C = 1)$$

$$p(Y) = q(Y|C = 0)$$

$$q(Y) = q(Y|C = 0)q(C = 0) + q(Y|C = 1)q(C = 1)$$

Density Ratio Trick:

$$\frac{p(Y|X)}{p(Y)} = \frac{q(Y|C = 1)}{q(Y|C = 0)} = \frac{\frac{q(C=1|Y)q(Y)}{q(C=1)}}{\frac{q(C=0|Y)q(Y)}{q(C=0)}} = \frac{q(C = 1|Y)}{q(C = 0|Y)}$$

Binary Classifier!

# Deep Graph InfoMax

How to estimate Mutual Information?

Suppose we generate  $Y$  from a mixture distribution:

$$q(C = 1) = q(C = 0) = \frac{1}{2}$$

$$p(Y|X) = q(Y|C = 1)$$

$$p(Y) = q(Y|C = 0)$$

$$q(Y) = q(Y|C = 0)q(C = 0) + q(Y|C = 1)q(C = 1)$$

Density Ratio Trick:

$$\frac{p(Y|X)}{p(Y)} = \frac{q(Y|C = 1)}{q(Y|C = 0)} = \frac{\frac{q(C=1|Y)q(Y)}{q(C=1)}}{\frac{q(C=0|Y)q(Y)}{q(C=0)}} = \frac{q(C = 1|Y)}{q(C = 0|Y)}$$

Binary Classifier!

We only need samplers to train a classifier, no exact probability densities are required!

# Deep Graph InfoMax

Given samples  $Y, C \sim q(Y, C)$

$$q(C = 1) = q(C = 0) = \frac{1}{2}$$

$$p(Y|X) = q(Y|C = 1)$$

$$p(Y) = q(Y|C = 0)$$

We build a classifier  $q(C = 1|Y) \propto f_{\theta}(X, Y)$

# Deep Graph InfoMax

Given samples  $Y, C \sim q(Y, C)$

$$q(C = 1) = q(C = 0) = \frac{1}{2}$$

$$p(Y|X) = q(Y|C = 1)$$

$$p(Y) = q(Y|C = 0)$$

We build a classifier  $q(C = 1|Y) \propto f_{\theta}(X, Y)$

Minimize Binary Cross-entropy:  $\min_q -\frac{1}{N} \left( \sum_{Y,C} C \log q(C = 1|Y) + (1 - C) \log(1 - q(C = 1|Y)) \right)$

# Deep Graph InfoMax

Given samples  $Y, C \sim q(Y, C)$

$$q(C = 1) = q(C = 0) = \frac{1}{2}$$

$$p(Y|X) = q(Y|C = 1)$$

$$p(Y) = q(Y|C = 0)$$

We build a classifier  $q(C = 1|Y) \propto f_{\theta}(X, Y)$

Minimize Binary Cross-entropy:  $\min_q -\frac{1}{N} \left( \sum_{Y,C} C \log q(C = 1|Y) + (1 - C) \log(1 - q(C = 1|Y)) \right)$

The optimal  $q$  is:

$$q^*(C = 1|Y) = \frac{\frac{p(Y|X)}{p(Y)}}{1 + \frac{p(Y|X)}{p(Y)}}$$

**Optimal Bayesian Classifier!**



# Deep Graph InfoMax

The optimal  $q$  is

$$q^*(C = 1|Y) = \frac{\frac{p(Y|X)}{p(Y)}}{1 + \frac{p(Y|X)}{p(Y)}}$$

Binary Cross-entropy:  $\mathcal{L}_{\text{BCE}}(q) = -\mathbb{E}_{C,Y} [C \log q(C = 1|Y) + (1 - C) \log(1 - q(C = 1|Y))]$

# Deep Graph InfoMax

The optimal  $q$  is

$$q^*(C = 1|Y) = \frac{\frac{p(Y|X)}{p(Y)}}{1 + \frac{p(Y|X)}{p(Y)}}$$

Binary Cross-entropy:  $\mathcal{L}_{\text{BCE}}(q) = -\mathbb{E}_{C,Y} [C \log q(C = 1|Y) + (1 - C) \log(1 - q(C = 1|Y))]$

One can show [4] that  $\mathcal{L}_{\text{BCE}}(q^*) \geq -\text{MI}(X, Y)$

One can also generalize this lower bound to multiple variables [4,5]

# Deep Graph InfoMax

We want to learn node representations that capture global context information of the graph, i.e., maximize *local mutual information*

$$\mathcal{L} = \frac{1}{N + M} \left( \sum_{i=1}^N \mathbb{E}_{(\mathbf{X}, \mathbf{A})} \left[ \log \mathcal{D} \left( \vec{h}_i, \vec{s} \right) \right] + \sum_{j=1}^M \mathbb{E}_{(\tilde{\mathbf{X}}, \tilde{\mathbf{A}})} \left[ \log \left( 1 - \mathcal{D} \left( \vec{h}_j, \vec{s} \right) \right) \right] \right)$$

- Graph input

# Deep Graph InfoMax

We want to learn node representations that capture global context information of the graph, i.e., maximize *local mutual information*

$$\mathcal{L} = \frac{1}{N + M} \left( \sum_{i=1}^N \mathbb{E}_{(\mathbf{X}, \mathbf{A})} \left[ \log \mathcal{D} \left( \vec{h}_i, \vec{s} \right) \right] + \sum_{j=1}^M \mathbb{E}_{(\tilde{\mathbf{X}}, \tilde{\mathbf{A}})} \left[ \log \left( 1 - \mathcal{D} \left( \vec{h}_j, \vec{s} \right) \right) \right] \right)$$

- Graph input
- Corrupted/Noisy graph input

# Deep Graph InfoMax

We want to learn node representations that capture global context information of the graph, i.e., maximize *local mutual information*

$$\mathcal{L} = \frac{1}{N + M} \left( \sum_{i=1}^N \mathbb{E}_{(\mathbf{X}, \mathbf{A})} \left[ \log \mathcal{D} \left( \vec{h}_i, \vec{s} \right) \right] + \sum_{j=1}^M \mathbb{E}_{(\tilde{\mathbf{X}}, \tilde{\mathbf{A}})} \left[ \log \left( 1 - \mathcal{D} \left( \vec{h}_j, \vec{s} \right) \right) \right] \right)$$

- Graph input
- Corrupted/Noisy graph input
- Node representation

# Deep Graph InfoMax

We want to learn node representations that capture global context information of the graph, i.e., maximize *local mutual information*

$$\mathcal{L} = \frac{1}{N + M} \left( \sum_{i=1}^N \mathbb{E}_{(\mathbf{X}, \mathbf{A})} \left[ \log \mathcal{D} \left( \vec{h}_i, \vec{s} \right) \right] + \sum_{j=1}^M \mathbb{E}_{(\tilde{\mathbf{X}}, \tilde{\mathbf{A}})} \left[ \log \left( 1 - \mathcal{D} \left( \vec{h}_j, \vec{s} \right) \right) \right] \right)$$

- Graph input
- Corrupted/Noisy graph input
- Node representation
- **Graph representation**

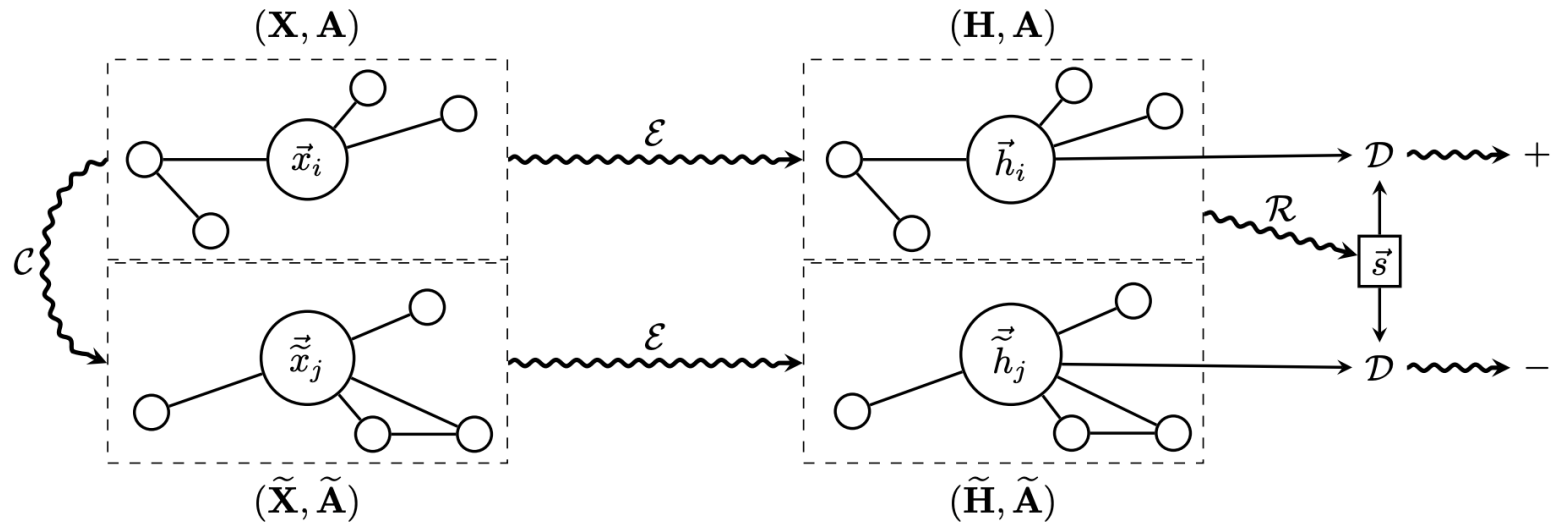
# Deep Graph InfoMax

We want to learn node representations that capture global context information of the graph, i.e., maximize *local mutual information*

$$\mathcal{L} = \frac{1}{N + M} \left( \sum_{i=1}^N \mathbb{E}_{(\mathbf{X}, \mathbf{A})} \left[ \log \mathcal{D}(\vec{h}_i, \vec{s}) \right] + \sum_{j=1}^M \mathbb{E}_{(\tilde{\mathbf{X}}, \tilde{\mathbf{A}})} \left[ \log \left( 1 - \mathcal{D}(\vec{h}_j, \vec{s}) \right) \right] \right)$$

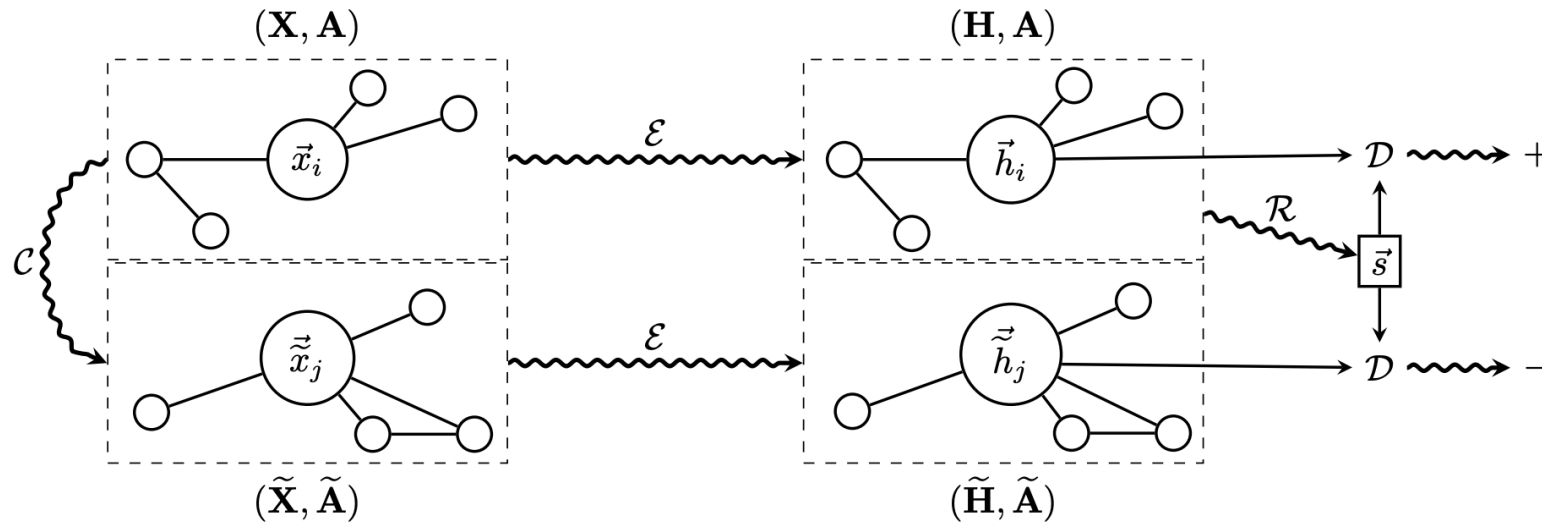
- Graph input
- Corrupted/Noisy graph input
- Node representation
- Graph representation
- **Discriminator / Binary classifier**

# Deep Graph InfoMax





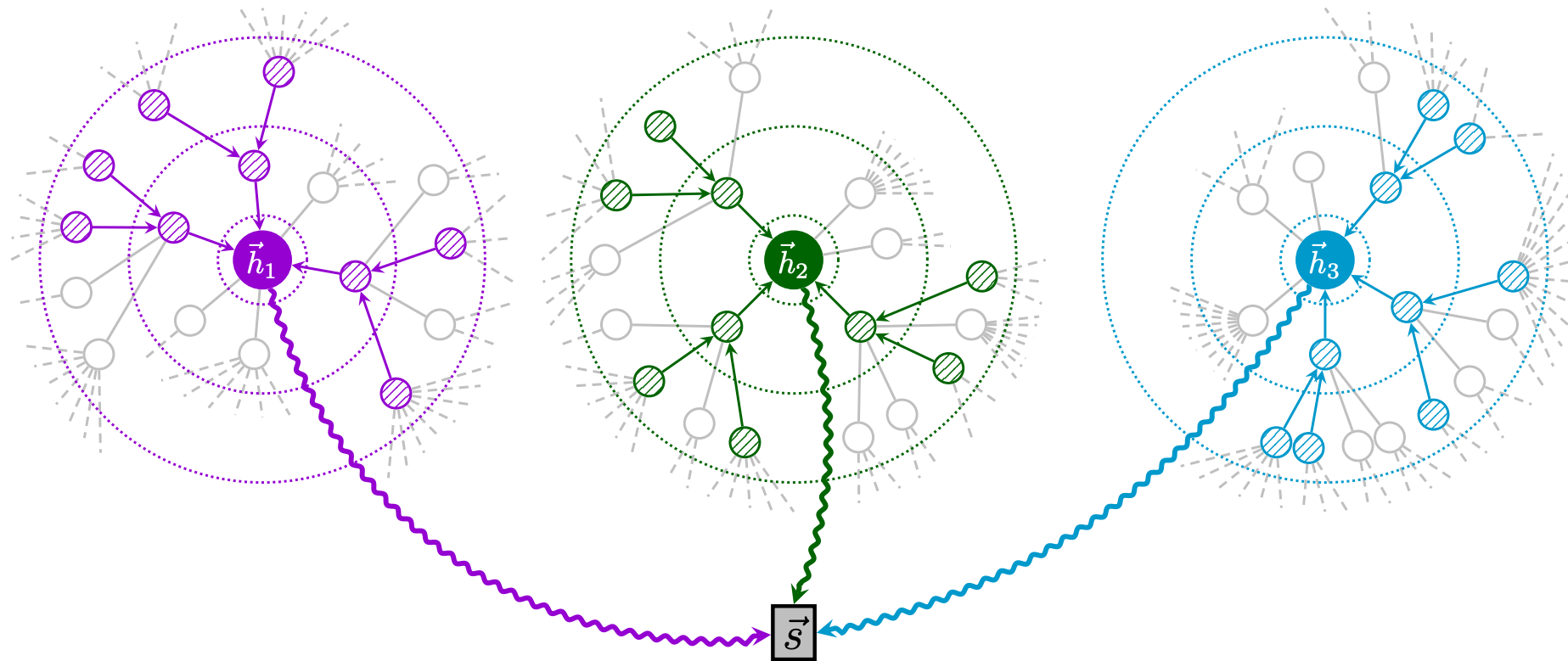
# Deep Graph InfoMax



1. Sample a negative example by using the corruption function:  $(\tilde{\mathbf{X}}, \tilde{\mathbf{A}}) \sim \mathcal{C}(\mathbf{X}, \mathbf{A})$ .
2. Obtain patch representations,  $\vec{h}_i$  for the input graph by passing it through the encoder:  $\mathbf{H} = \mathcal{E}(\mathbf{X}, \mathbf{A}) = \{\vec{h}_1, \vec{h}_2, \dots, \vec{h}_N\}$ .
3. Obtain patch representations,  $\vec{h}_j$  for the negative example by passing it through the encoder:  $\tilde{\mathbf{H}} = \mathcal{E}(\tilde{\mathbf{X}}, \tilde{\mathbf{A}}) = \{\vec{h}_1, \vec{h}_2, \dots, \vec{h}_M\}$ .
4. Summarize the input graph by passing its patch representations through the readout function:  $\vec{s} = \mathcal{R}(\mathbf{H})$ .
5. Update parameters of  $\mathcal{E}$ ,  $\mathcal{R}$  and  $\mathcal{D}$  by applying gradient descent to maximize Equation 1.

# Deep Graph InfoMax

If a huge graph is presented, sampling subgraphs is necessary



# Deep Graph InfoMax

<i>Transductive</i>				
Available data	Method	Cora	Citeseer	Pubmed
X	Raw features	47.9 ± 0.4%	49.3 ± 0.2%	69.1 ± 0.3%
A, Y	LP (Zhu et al., 2003)	68.0%	45.3%	63.0%
A	DeepWalk (Perozzi et al., 2014)	67.2%	43.2%	65.3%
X, A	DeepWalk + features	70.7 ± 0.6%	51.4 ± 0.5%	74.3 ± 0.9%
X, A	Random-Init (ours)	69.3 ± 1.4%	61.9 ± 1.6%	69.6 ± 1.9%
X, A	<b>DGI</b> (ours)	<b>82.3 ± 0.6%</b>	<b>71.8 ± 0.7%</b>	<b>76.8 ± 0.6%</b>
X, A, Y	GCN (Kipf & Welling, 2016a)	81.5%	70.3%	79.0%
X, A, Y	Planetoid (Yang et al., 2016)	75.7%	64.7%	77.2%

<i>Inductive</i>			
Available data	Method	Reddit	PPI
X	Raw features	0.585	0.422
A	DeepWalk (Perozzi et al., 2014)	0.324	—
X, A	DeepWalk + features	0.691	—
X, A	GraphSAGE-GCN (Hamilton et al., 2017a)	0.908	0.465
X, A	GraphSAGE-mean (Hamilton et al., 2017a)	0.897	0.486
X, A	GraphSAGE-LSTM (Hamilton et al., 2017a)	0.907	0.482
X, A	GraphSAGE-pool (Hamilton et al., 2017a)	0.892	0.502
X, A	Random-Init (ours)	0.933 ± 0.001	0.626 ± 0.002
X, A	<b>DGI</b> (ours)	<b>0.940 ± 0.001</b>	<b>0.638 ± 0.002</b>
X, A, Y	FastGCN (Chen et al., 2018)	0.937	—
X, A, Y	Avg. pooling (Zhang et al., 2018)	0.958 ± 0.001	0.969 ± 0.002

# SimCLR

# SimCLR

---

**Algorithm 1** SimCLR's main learning algorithm.

---

**input:** batch size  $N$ , constant  $\tau$ , structure of  $f$ ,  $g$ ,  $\mathcal{T}$ .  
**for** sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  **do**  
  **for all**  $k \in \{1, \dots, N\}$  **do**  
    draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$   
    # the first augmentation  
     $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$   
     $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation  
     $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection  
    # the second augmentation  
     $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$   
     $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation  
     $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection  
  **end for**  
  **for all**  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  **do**  
     $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity  
  **end for**  
  **define**  $\ell(i, j)$  **as**  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbf{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$   
   $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$   
  update networks  $f$  and  $g$  to minimize  $\mathcal{L}$   
**end for**  
**return** encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$

---

# SimCLR



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate  $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise

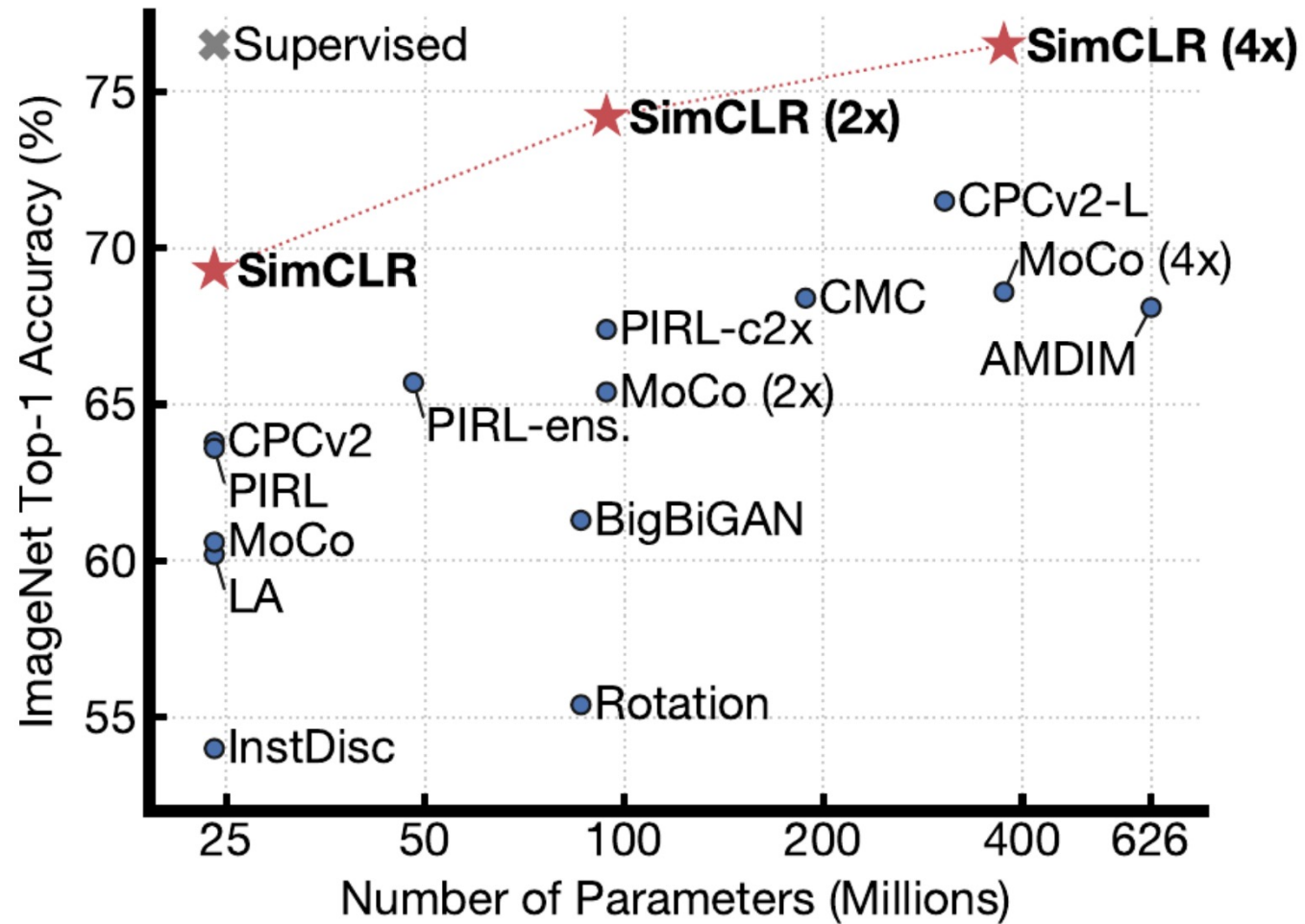


(i) Gaussian blur



(j) Sobel filtering

# SimCLR



# References

- [1] Perozzi, B., Al-Rfou, R. and Skiena, S., 2014, August. Deepwalk: Online learning of social representations. In Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining (pp. 701-710).
- [2] Grover, A. and Leskovec, J., 2016, August. node2vec: Scalable feature learning for networks. In Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining (pp. 855-864).
- [3] Veličković, P., Fedus, W., Hamilton, W.L., Liò, P., Bengio, Y. and Hjelm, R.D., 2018. Deep graph infomax. arXiv preprint arXiv:1809.10341.
- [4] Van den Oord, A., Li, Y. and Vinyals, O., 2018. Representation learning with contrastive predictive coding. arXiv e-prints, pp.arXiv-1807.
- [5] Hjelm, R.D., Fedorov, A., Lavoie-Marchildon, S., Grewal, K., Bachman, P., Trischler, A. and Bengio, Y., 2018. Learning deep representations by mutual information estimation and maximization. arXiv preprint arXiv:1808.06670.
- [6] Chen T, Kornblith S, Norouzi M, Hinton G. A simple framework for contrastive learning of visual representations. In International conference on machine learning 2020 Nov 21 (pp. 1597-1607). PMLR.



Questions?