

# DGI

## Deep Graph Infomax

이수연

2022.02.03

# Contents

---

1. Backgrounds
2. DGI Methodology
3. Classification Performance
4. Conclusion

# Backgrounds

the dominant algorithms for unsupervised representation learning with graph-structured data rely on **random walk**-based objectives

nodes that are “close” in the input graph are also “close” in the representation space

## Limitation

1. over-emphasizing **proximity** information at the expense of **structural** information
2. highly dependent on hyperparameter
3. with stronger encoder, unclear whether it actually provide any useful signal

→ Alternative objective based on **mutual information**

# Backgrounds

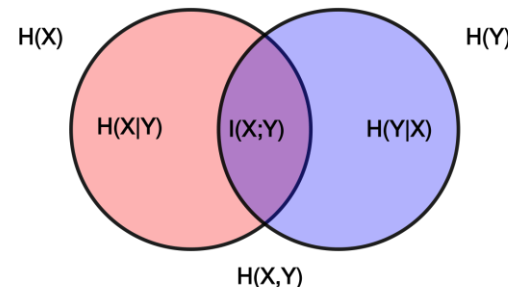
## Mutual Information(MI)

measures how much one random variables tells us about another (dependence)

$$I(X; Y) = \sum_{x,y} P_{XY}(x, y) \log \frac{P_{XY}(x, y)}{P_X(x)P_Y(y)} = E_{P_{XY}} \log \frac{P_{XY}}{P_X P_Y}$$

distribution을 알기 어려워서 MI 직접 계산 힘들

→ lower-bound를 maximize,  $I(X, Z) \geq I_{\theta}(X, Z)$



**DV** (Donsker-Varadhan)  $\mathcal{I}(X; Y) := \mathcal{D}_{KL}(\mathbb{J}||\mathbb{M}) \geq \hat{\mathcal{I}}_{\omega}^{(DV)}(X; Y) := \mathbb{E}_{\mathbb{J}}[T_{\omega}(x, y)] - \log \mathbb{E}_{\mathbb{M}}[e^{T_{\omega}(x, y)}]$  **MINE**  
(Mutual Information Neural Estimation)

**JSD** (Jensen-Shannon)  $\hat{\mathcal{I}}_{\omega, \psi}^{(JSD)}(X; E_{\psi}(X)) := \mathbb{E}_{\mathbb{P}}[-\text{sp}(-T_{\psi, \omega}(x, E_{\psi}(x)))] - \mathbb{E}_{\mathbb{P} \times \tilde{\mathbb{P}}}[\text{sp}(T_{\psi, \omega}(x', E_{\psi}(x)))]$

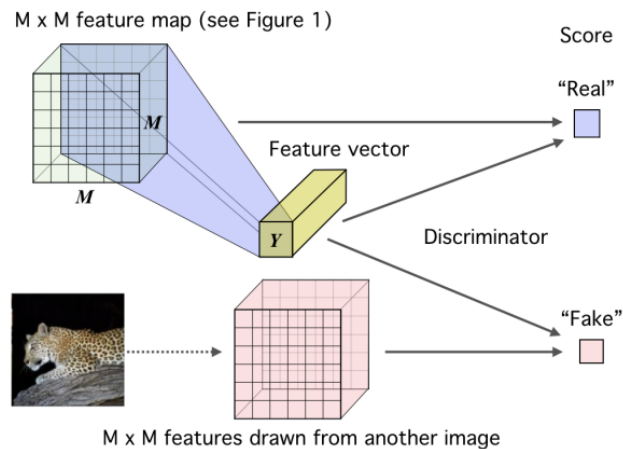
**infoNCE** (Noise-Contrastive Estimation + softmax)  $\hat{\mathcal{I}}_{\omega, \psi}^{(\text{infoNCE})}(X; E_{\psi}(X)) := \mathbb{E}_{\mathbb{P}} \left[ T_{\psi, \omega}(x, E_{\psi}(x)) - \mathbb{E}_{\tilde{\mathbb{P}}} \left[ \log \sum_{x'} e^{T_{\psi, \omega}(x', E_{\psi}(x))} \right] \right]$

# Backgrounds

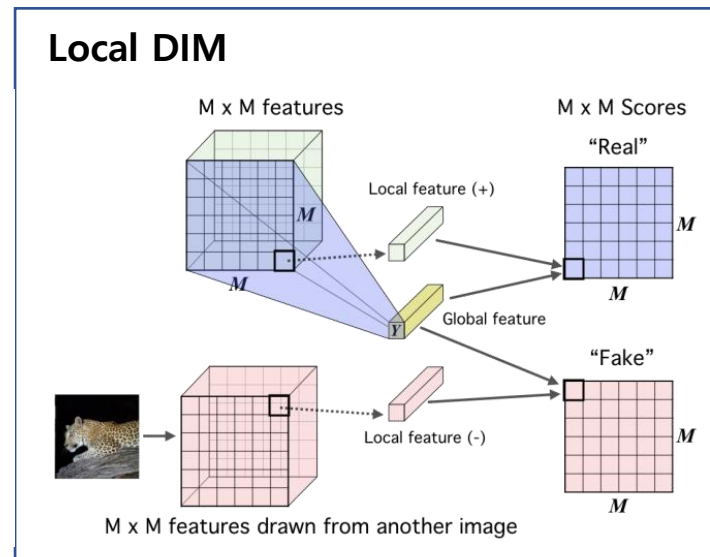
## Deep InfoMax (DIM)

unsupervised learning of representations by maximizing mutual information between an input and the output of a deep neural network encoder

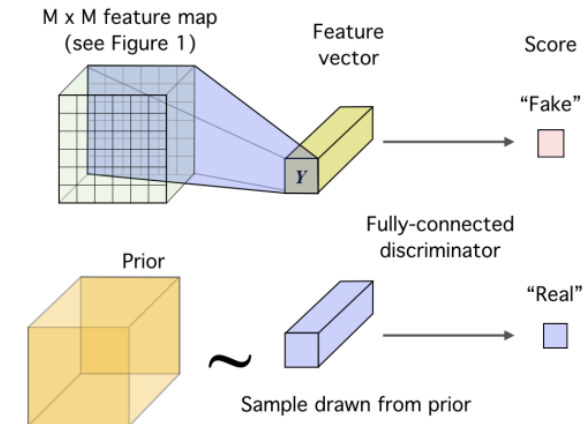
### Global DIM



### Local DIM



### matching to prior

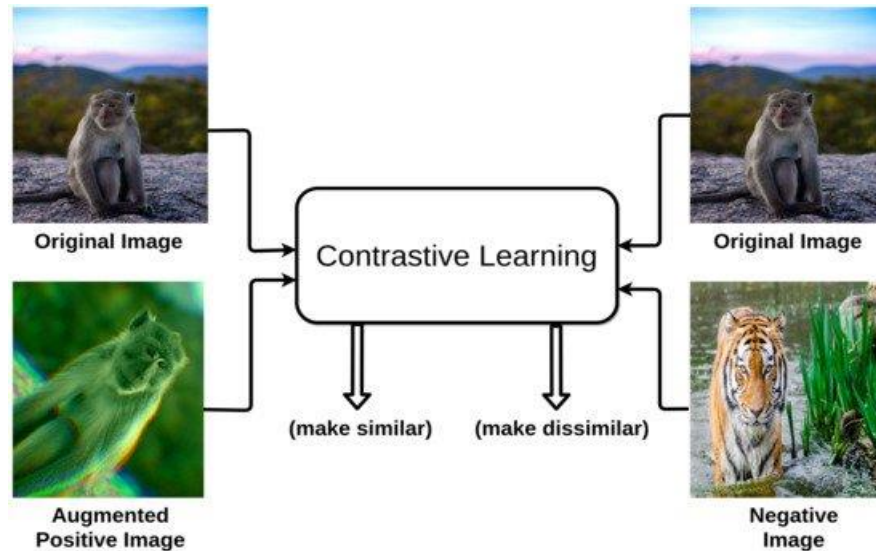


**objective** 
$$\arg \max_{\omega_1, \omega_2, \psi} (\alpha \hat{\mathcal{I}}_{\omega_1, \psi}(X; E_{\psi}(X)) + \frac{\beta}{M^2} \sum_{i=1}^{M^2} \hat{\mathcal{I}}_{\omega_2, \psi}(X^{(i)}; E_{\psi}(X))) + \arg \min_{\psi} \arg \max_{\phi} \gamma \hat{\mathcal{D}}_{\phi}(\mathbb{V} || \mathbb{U}_{\psi, \mathbb{P}})$$

# Backgrounds

## Contrastive Learning

an important approach for **self-supervised** learning of representations



Basic intuition behind contrastive learning paradigm  
from: A Survey on Contrastive Self-Supervised Learning

"push original and augmented images closer  
and push original and negative images away"

real(i.e. positive)과 fake(i.e. negative) example  
을 대조(**contrast**)하여 representation 학습

# Backgrounds

---

**DIM을 그래프에 적용!**

- local features to capture information shared across the entire graph
- MI maximization
- Contrastive learning

# DGI Methodology

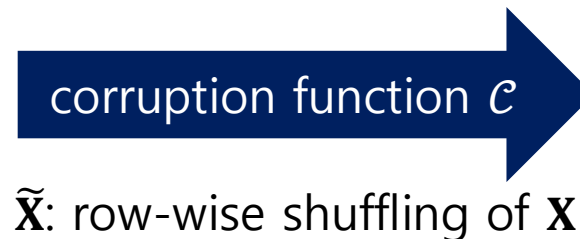
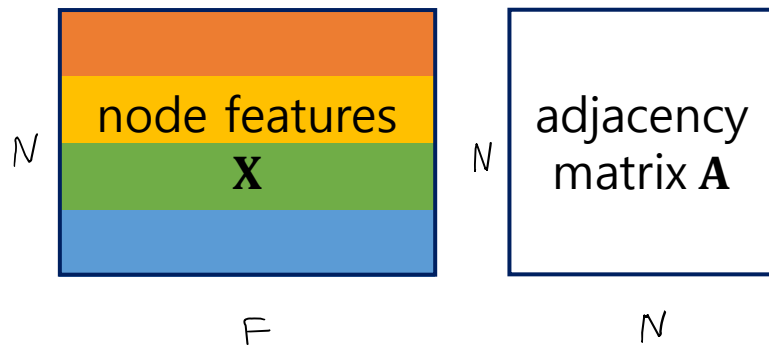
## Negative Sampling

### multi-graph setting

corruption function simply samples a different graph from the training set

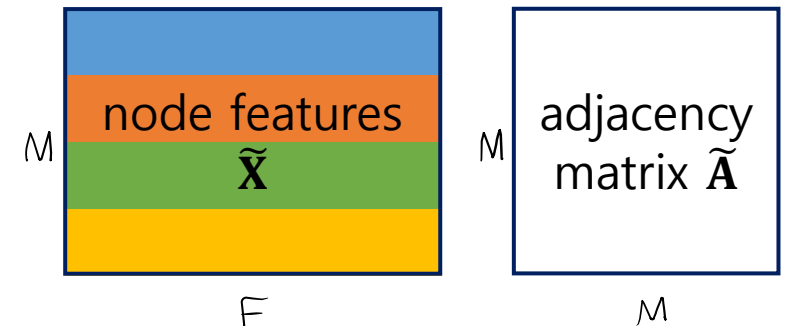
### single-graph setting

original graph  $(\mathbf{X}, \mathbf{A})$



$\tilde{\mathbf{X}}$ : row-wise shuffling of  $\mathbf{X}$

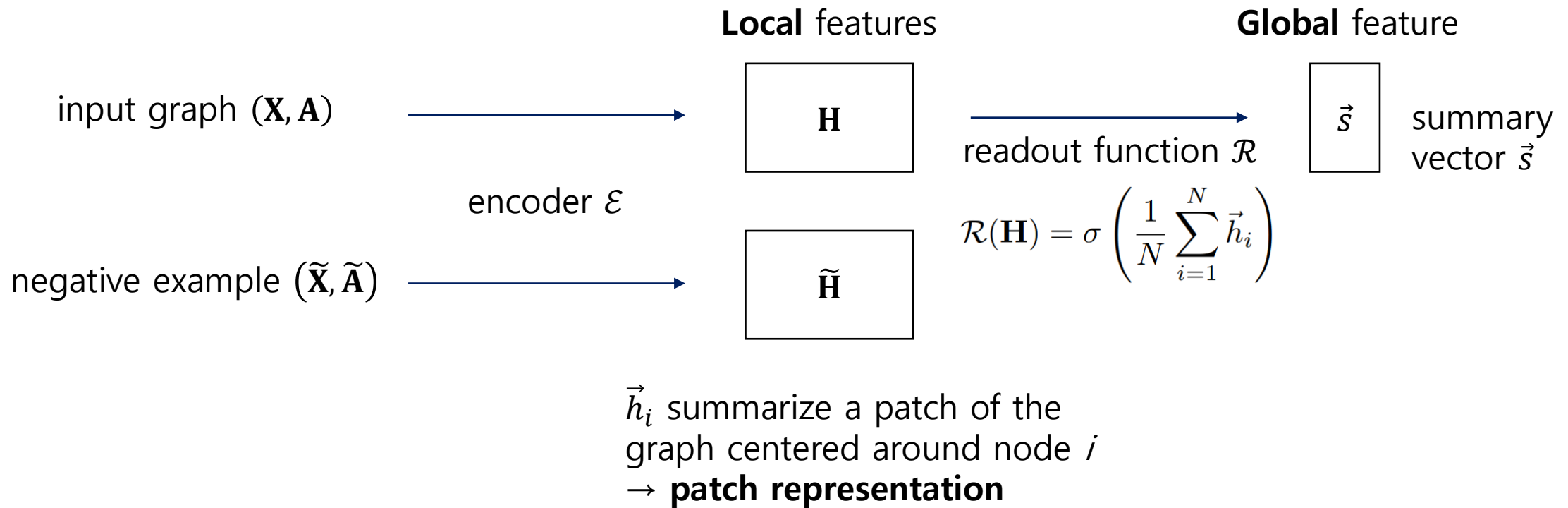
alternative graph  $(\tilde{\mathbf{X}}, \tilde{\mathbf{A}})$





# DGI Methodology

## Local & Global features



# DGI Methodology

## Theoretical Motivation

**Theorem 1.** minimizing the classification error in the discriminator  $\equiv$  maximizing the mutual information between the input and output

$$\vec{s}^* = \operatorname{argmin}_{\vec{s}} \operatorname{Err}^*$$

$$\vec{s}^* = \operatorname{argmax}_{\vec{s}} \operatorname{MI}(\mathbf{X}; \vec{s})$$

**Theorem 2.** the  $\vec{h}_i$  that minimizes the classification error between  $p(\vec{h}_i, \vec{s})$  and  $p(\vec{h}_i)p(\vec{s})$  also maximizes  $\operatorname{MI}(\mathbf{X}_i^{(k)}; \vec{h}_i)$

→ a classifier between samples from the joint(positive examples) and the product of marginals (negative examples), and using the binary cross-entropy (BCE) loss to optimize this classifier

# DGI Methodology

## Mutual Information Maximization

discriminator  $\mathcal{D}(\vec{h}_i, \vec{s}) = \sigma(\vec{h}_i^T \mathbf{W} \vec{s})$

positive example과 negative example을 구별하도록 training

Glorot initialization  
Adam SGD optimizer

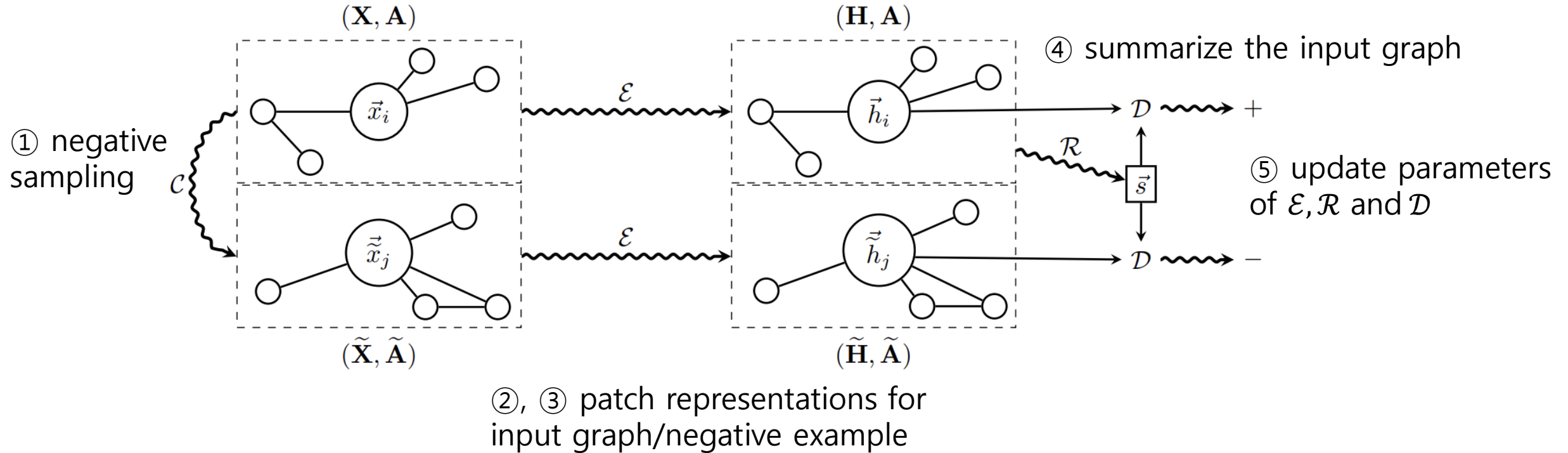
$$\mathcal{L} = \frac{1}{N + M} \left( \sum_{i=1}^N \mathbb{E}_{(\mathbf{X}, \mathbf{A})} \left[ \log \mathcal{D}(\vec{h}_i, \vec{s}) \right] + \sum_{j=1}^M \mathbb{E}_{(\tilde{\mathbf{X}}, \tilde{\mathbf{A}})} \left[ \log \left( 1 - \mathcal{D}(\vec{\tilde{h}}_j, \vec{s}) \right) \right] \right)$$

noise-contrastive type objective with a standard binary cross-entropy (BCE) loss

→ mutual information with the global graph summary, patch-level의 similarity

# DGI Methodology

## Overview of DGI



# Classification Performance

patch representations learned in a fully unsupervised manner  
→ evaluating the node-level classification utility

## Datasets

- (1) Cora, Citeseer, Pubmed citation networks: classifying research papers into topics (transductive)
- (2) Reddit social network: predicting the community structure (inductive)
- (3) PPI network: classifying protein roles (inductive, multiple graphs)

# Classification Performance

## Encoder

### Transductive learning

one-layer GCN

$$\mathcal{E}(\mathbf{X}, \mathbf{A}) = \sigma \left( \hat{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X} \Theta \right)$$

the learned filters rely on a fixed and known adjacency matrix: not appropriate to Inductive learning

### Inductive learning on large graphs

three-layer mean-pooling model with skip connections

$$\text{MP}(\mathbf{X}, \mathbf{A}) = \hat{\mathbf{D}}^{-1} \hat{\mathbf{A}} \mathbf{X} \Theta$$

$$\mathcal{E}(\mathbf{X}, \mathbf{A}) = \widetilde{\text{MP}}_3(\widetilde{\text{MP}}_2(\widetilde{\text{MP}}_1(\mathbf{X}, \mathbf{A}), \mathbf{A}), \mathbf{A})$$

$\mathbf{X}$ ,  $\mathbf{A}$  are not needed

sampling node neighborhoods with replacement

### Inductive learning on multiple graphs

three-layer mean-pooling model with dense skip connections

$$\mathbf{H}_1 = \sigma(\text{MP}_1(\mathbf{X}, \mathbf{A}))$$

$$\mathbf{H}_2 = \sigma(\text{MP}_2(\mathbf{H}_1 + \mathbf{X} \mathbf{W}_{\text{skip}}, \mathbf{A}))$$

$$\mathcal{E}(\mathbf{X}, \mathbf{A}) = \sigma(\text{MP}_3(\mathbf{H}_2 + \mathbf{H}_1 + \mathbf{X} \mathbf{W}_{\text{skip}}, \mathbf{A}))$$

# Classification Performance

## Result

the DGI approach is competitive with the results reported for the GCN model with the supervised loss

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

note that... supervised transductive SOTA(GraphSGAN)는 못 넘음

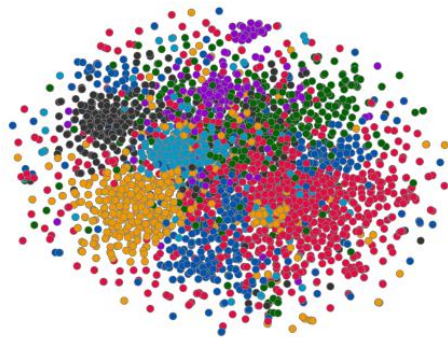
		<i>Inductive</i> micro-averaged F <sub>1</sub> score	
Available data	Method	Reddit	PPI
<b>X</b>	Raw features	0.585	0.422
<b>A</b>	DeepWalk (Perozzi et al., 2014)	0.324	—
<b>X, A</b>	DeepWalk + features	0.691	—
<b>X, A</b>	GraphSAGE-GCN (Hamilton et al., 2017a)	0.908	0.465
<b>X, A</b>	GraphSAGE-mean (Hamilton et al., 2017a)	0.897	0.486
<b>X, A</b>	GraphSAGE-LSTM (Hamilton et al., 2017a)	0.907	0.482
<b>X, A</b>	GraphSAGE-pool (Hamilton et al., 2017a)	0.892	0.502
<b>X, A</b>	Random-Init (ours)	0.933 ± 0.001	0.626 ± 0.002
<b>X, A</b>	<b>DGI (ours)</b>	<b>0.940 ± 0.001</b>	<b>0.638 ± 0.002</b>
<b>X, A, Y</b>	FastGCN (Chen et al., 2018)	0.937	—
<b>X, A, Y</b>	Avg. pooling (Zhang et al., 2018)	0.958 ± 0.001	0.969 ± 0.002

potential in the inductive node classification

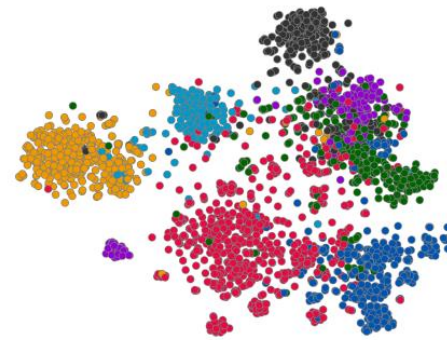
large gap... attributed to the extreme sparsity of available node features

# Classification Performance

t-SNE embeddings of the nodes in the Cora dataset



raw features



learned DGI model

Silhouette score of 0.234



# Conclusions

---

**Deep InfoMax** → Deep Graph Infomax

**local mutual information maximization** across the graph's patch representations  
→ node embeddings that are mindful of the global structural properties of the graph

competitive performance across a variety of both **transductive** and **inductive** classification tasks

감사합니다