

Introduction to GNNs and Self-supervised Learning on Graphs

I. Makarov & V. Pozdnyakov & D. Kiselev

BigData Academy MADE from Mail.ru Group

Graph Neural Networks and Applications



- Instructor: Ilya Makarov
- Tutor: Vitaly Pozdnyakov, Dmitrii Kiselev
- Invited lecturers: to be announced
- Course length: 11-12 lectures/classes
- Lecturer's Telegram: @iamakarov (urgent)
- Tutor's Telegram: @pozdnyakov_vitaliy (urgent), @dkiselev (urgent)
- Discord (questions, deadlines, grading)
- Programming: Python, iPython notebooks, Anaconda distribution
- Python libraries: NetworkX, pyG, DGL

Prerequisites

- Network Science
- Deep Learning
- Machine Learning
- Differential Equations
- Programming in Python

- "CS224W: Machine Learning with Graphs". Jure Leskovec. Stanford, 2020-2021.
- "SNAP Project". <http://snap.stanford.edu/>
- "PyG & PyG Temporal".
<https://pytorch-geometric.readthedocs.io/en/latest/>,
<https://pytorch-geometric-temporal.readthedocs.io/en/latest/>
- "Network Science", Albert-Laszlo Barabasi, Cambridge University Press, 2016. <http://networksciencebook.com>
- "Network Science" course by Leonid E. Zhukov.

Topics (max list)

- 1 Self-supervised learning on graphs
- 2 Subgraph Embedding & Deep Sets
- 3 Scaling GNNs: Efficient models, Sampling models
- 4 Compressing GNNs: Quantization, Distillation
- 5 Deep Generative Graph Models & adversarial robustness
- 6 Temporal graph embeddings
- 7 GNN Explanations
- 8 GNN & RecSys & KG
- 9 Query Embeddings for KGs
- 10 Transport GNNs + Combinatorics
- 11 Guest lectures: Applications to Biology, Medicine, Bioinformatics, Chemistry, Physics
- 12 OGB benchmark (optional challenge)

Recap

» Go to Lecture 1: SSL on Graphs

- ① Real-world networks
- ② Random graphs
- ③ Centrality & Prestige node measures
- ④ Structural similarity in networks
- ⑤ Community detection
- ⑥ Modeling dynamics on graphs
 - continuous: diffusion, epidemics
 - discrete: information propagation, segregation
- ⑦ Graph Machine Learning:
 - label propagation and node similarities
 - node and edge embeddings
 - graph neural networks
- ⑧ Knowledge Graphs

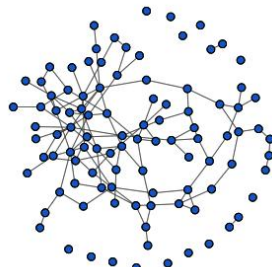
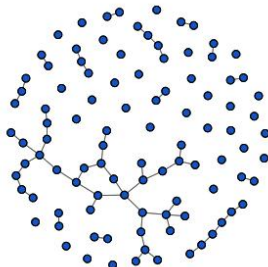
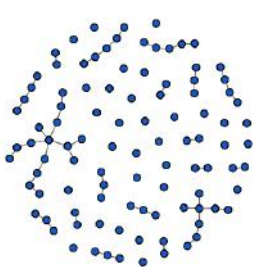
Complex networks (nor regular, nor random)

- ① Power law node degree distribution: "scale-free" networks
- ② Small diameter and average path length: "small world" networks
- ③ High clustering coefficient: transitivity

Random graph model

Consider $G_{n,p}$ as a function of p

- $p = 0$, empty graph - $\langle k \rangle = 0$
- $p = 1$, complete (full) graph - $\langle k \rangle = n - 1$
- n_G -largest connected component, $s = \frac{n_G}{n}$

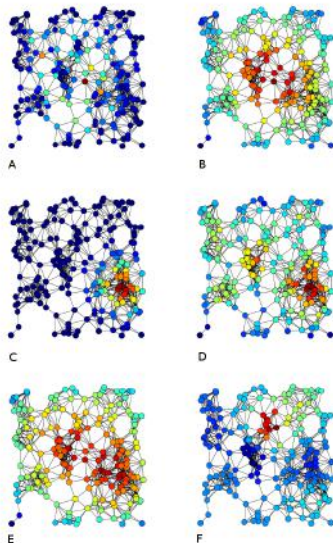


p

Random graph models comparison

	Random	BA model	WS model	Empirical networks
$P(k)$	$\frac{\lambda^k e^{-\lambda}}{k!}$	k^{-3}	poisson like	power law
C	$\langle k \rangle / N$	$N^{-0.75}$	const	large
$\langle L \rangle$	$\frac{\log(N)}{\log(\langle k \rangle)}$	$\frac{\log(N)}{\log \log(N)}$	$\log(N)$	small

Centrality examples

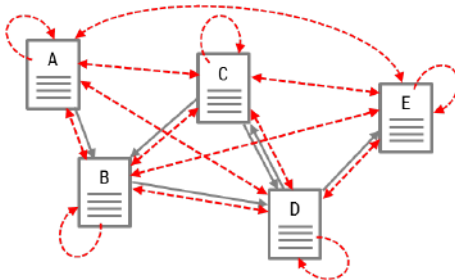


from Wikipedia

- A) Betweenness centrality
- B) Closeness centrality
- C) Eigenvector centrality
- D) Degree centrality
- F) Harmonic centrality
- E) Katz centrality

PageRank

"PageRank can be thought of as a model of user behavior. We assume there is a "random surfer" who is given a web page at random and keeps clicking on links, never hitting "back" but eventually gets bored and starts on another random page. The **probability** that the random surfer visits a page is its **PageRank**."



Sergey Brin and Larry Page, 1998

- Power iterations:

$$\mathbf{p} \leftarrow \alpha \mathbf{P}^T \mathbf{p} + (1 - \alpha) \frac{\mathbf{e}}{n}, \quad \alpha - \text{teleportation coefficient}$$

- Sparse linear system:

$$(\mathbf{I} - \alpha \mathbf{P}^T) \mathbf{p} = (1 - \alpha) \frac{\mathbf{e}}{n}$$

- Eigenvalue problem ($\lambda = 1$):

$$\left(\alpha \mathbf{P}^T + (1 - \alpha) \mathbf{E} \right) \mathbf{p} = \lambda \mathbf{p}$$

$$\mathbf{P} = \mathbf{D}^{-1} \mathbf{A}$$

PageRank variations

- Power iterations

$$\mathbf{p} \leftarrow \alpha \mathbf{P}^T \mathbf{p} + (1 - \alpha) \mathbf{v}, \quad \mathbf{v} - \text{teleportation vector}$$

$$\mathbf{P}' = \alpha \mathbf{P} + (1 - \alpha) \mathbf{e} \mathbf{v}^T$$

$$\mathbf{p} \leftarrow \mathbf{P}'^T \mathbf{p}, \quad \|\mathbf{p}\| = 1$$

- Topic specific PageRank

\mathbf{v} - set of pages on specific topics

- TrustRank

\mathbf{v} - set of trusted pages

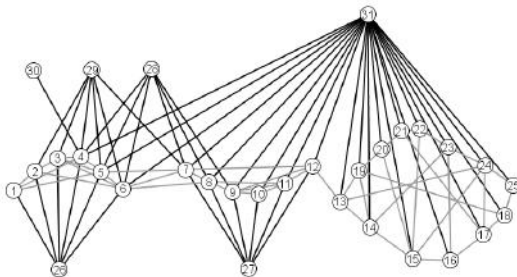
- Personalized PageRank

\mathbf{v} - set of personal preference pages

Structural similarity

Definition

Two nodes are similar to each other if they share many neighbors.



- Jaccard similarity

$$J(v_i, v_j) = \frac{|\mathcal{N}(v_i) \cap \mathcal{N}(v_j)|}{|\mathcal{N}(v_i) \cup \mathcal{N}(v_j)|}$$

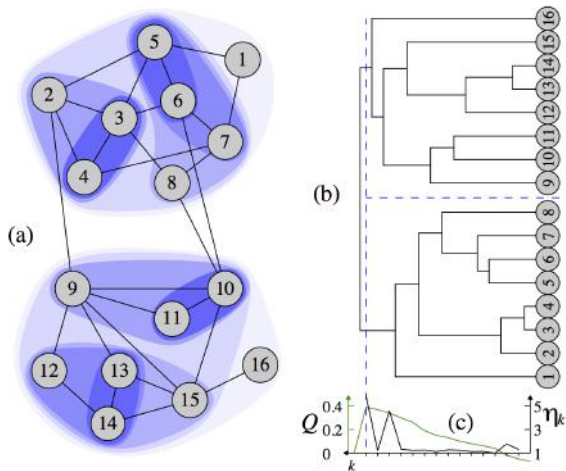
- Cosine similarity (vectors in n -dim space)

$$\sigma(v_i, v_j) = \cos(\theta_{ij}) = \frac{\mathbf{v}_i^T \mathbf{v}_j}{|\mathbf{v}_i| |\mathbf{v}_j|} = \frac{\sum_k A_{ik} A_{kj}}{\sqrt{\sum A_{ik}^2} \sqrt{\sum A_{jk}^2}}$$

- Pearson correlation coefficient:

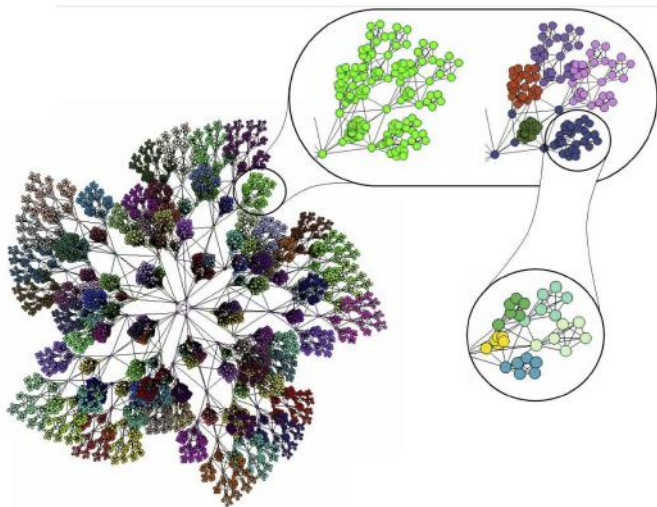
$$r_{ij} = \frac{\sum_k (A_{ik} - \langle A_i \rangle)(A_{jk} - \langle A_j \rangle)}{\sqrt{\sum_k (A_{ik} - \langle A_i \rangle)^2} \sqrt{\sum_k (A_{jk} - \langle A_j \rangle)^2}}$$

Walktrap



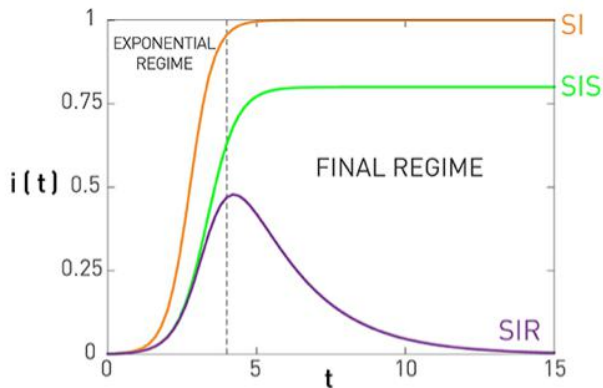
P. Pons and M. Latapy, 2006

Fast community unfolding



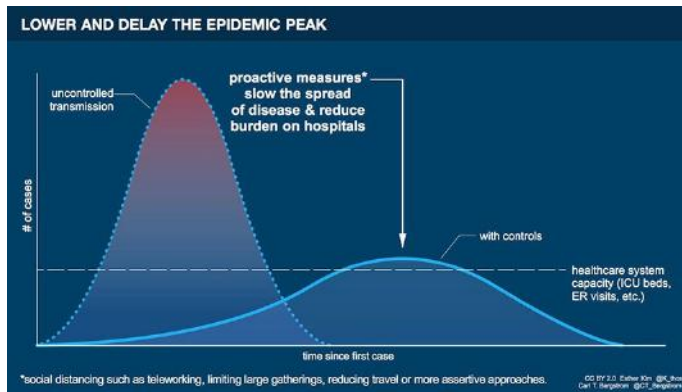
V. Blondel et.al., 2008

Compartmental models summary



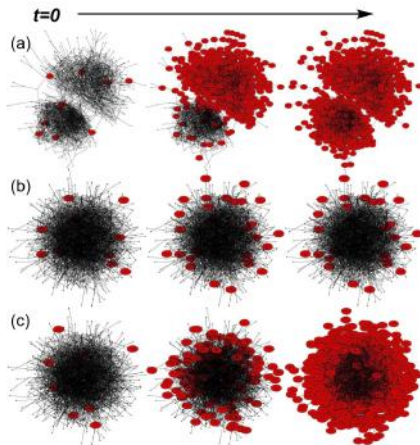
from Barabasi, 2016

Flatten the curve!



Cascades in random networks

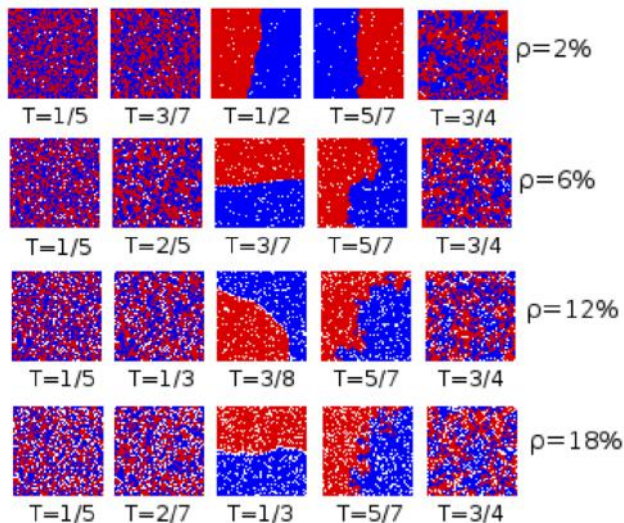
multiple seed nodes



(a) Empirical network; (b), (c) - randomized network

P. Singh, 2013

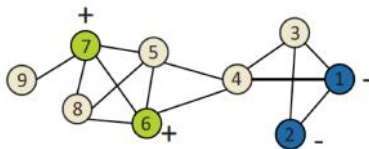
Spatial segregation



- Node classification (attribute inference)
- Link prediction (missing/hidden links inference)
- Community detection (clustering nodes in graph)
- Graph visualization (cluster projections)

Node classification

- Node classification - labeling of all nodes in a graph structure
- Subset of nodes is labeled: categorical/numeric/binary values
- Extend labeling to all nodes on the graph (class/class probability/regression)
- Classification in networked data, network classification, structured inference, relational learning



Label propagation

Algorithm: Label propagation, Zhu et. al 2002

Input: Graph $G(V, E)$, labels Y_I

Output: labels \hat{Y}

Compute $D_{ii} = \sum_j A_{ij}$

Compute $P = D^{-1}A$

Initialize $Y^{(0)} = (Y_I, 0)$, $t=0$

repeat

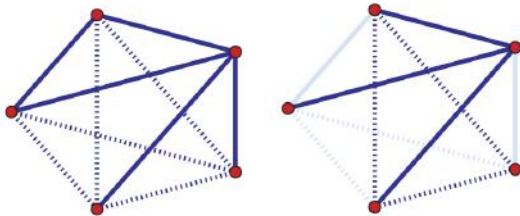
$Y^{(t+1)} \leftarrow P \cdot Y^{(t)}$
 $Y_I^{(t+1)} \leftarrow Y_I^{(t)}$

until $Y^{(t)}$ converges;

$\hat{Y} \leftarrow Y^{(t)}$

Solution: $\hat{Y} = \lim_{t \rightarrow \infty} Y^{(t)} = (I - P_{uu})^{-1} P_{ul} Y_I$

Link prediction



- Graph $G(V, E)$
- Number of "missing edges": $|V|(|V| - 1)/2 - |E|$
- In sparse graphs $|E| \ll |V|^2$, Prob. of correct random guess $O(\frac{1}{|V|^2})$

Local similarity indices

Local neighborhood of v_i and v_j

- Number of common neighbors:

$$s_{ij} = |\mathcal{N}(v_i) \cap \mathcal{N}(v_j)|$$

- Jaccard's coefficient:

$$s_{ij} = \frac{|\mathcal{N}(v_i) \cap \mathcal{N}(v_j)|}{|\mathcal{N}(v_i) \cup \mathcal{N}(v_j)|}$$

- Resource allocation:

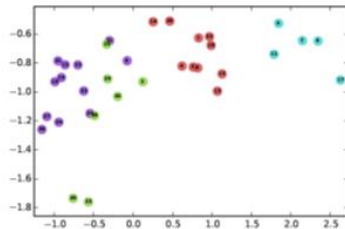
$$s_{ij} = \sum_{w \in \mathcal{N}(v_i) \cap \mathcal{N}(v_j)} \frac{1}{|\mathcal{N}(w)|}$$

Adamic/Adar:

$$s_{ij} = \sum_{w \in \mathcal{N}(v_i) \cap \mathcal{N}(v_j)} \frac{1}{\log |\mathcal{N}(w)|}$$

Graph Embeddings

- Necessity to automatically select features
- Reduce domain- and task- specific bias
- Unified framework to vectorize network
- Preserve graph properties in vector space
- Similar nodes \rightarrow close embeddings

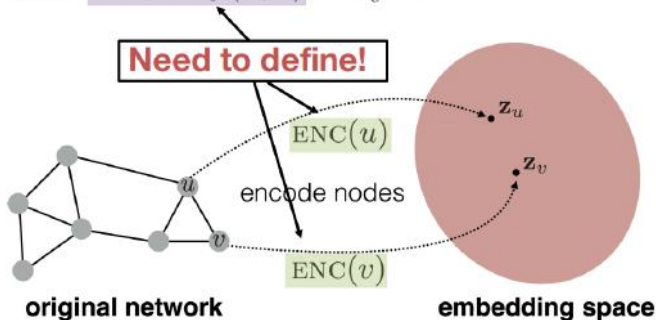


¹<http://snap.stanford.edu/proj/embeddings-www/>

Graph Embeddings

- Define **Encoder**
- Define **Similarity**/graph feature to preserve graph properties
- Define similarity/distance in the embedding space
- **Optimize** loss to fit embedding with similarity computed on graph

Goal: $\text{similarity}(u, v) \approx \mathbf{z}_v^\top \mathbf{z}_u$



First-order Proximity

- Similarity between u and v is A_{uv}
- MSE Loss
- Variant of Matrix Decomposition

The diagram illustrates the formula for the loss \mathcal{L} as a sum over all node pairs $(u, v) \in V \times V$ of the squared difference between the embedding similarity $\mathbf{z}_u^\top \mathbf{z}_v$ and the weighted adjacency matrix element $\mathbf{A}_{u,v}$. Colored boxes highlight the components: \mathcal{L} (red), the summation (green), the embedding similarity (purple), and the adjacency matrix element (blue). Arrows point from descriptive text to these components.

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \| \mathbf{z}_u^\top \mathbf{z}_v - \mathbf{A}_{u,v} \|^2$$

loss (what we want to minimize)

sum over all node pairs

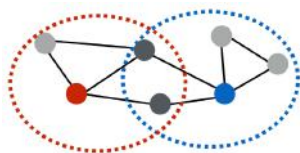
embedding similarity

(weighted) adjacency matrix for the graph

from Leskovec et al., 2018

Multi-order Proximity

- Similarity score S_{uv} as Jaccard/Common Neighbours, etc. (HOPE)



$$\mathcal{L} = \sum_{(u,v) \in V \times V} \left\| \boxed{\mathbf{z}_u^\top \mathbf{z}_v} - \boxed{S_{u,v}} \right\|^2$$

embedding similarity

multi-hop network similarity (i.e., any neighborhood overlap measure)

- Weighted k-hop paths with different k (GraRep)

$$\tilde{\mathbf{A}}_{i,j}^k = \max \left(\log \left(\frac{(\mathbf{A}_{i,j}/d_i)}{\sum_{l \in V} (\mathbf{A}_{l,j}/d_l)^k} \right)^k - \alpha, 0 \right)$$

node degree

constant shift

from Leskovec et al., 2018

- Even worse complexity

- Similarity between u and v is probability to co-occur on a random walk
- Sample each vertex u neighborhood $N_R(u)$ (multiset) by short random walks via strategy R
- Optimize similarity considering independent neighbor samples via MLE (remind Word2Vec)

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log(P(v|\mathbf{z}_u))$$

from Leskovec et al., 2018

- $P(v|z_u)$ is approximated via softmax over similarity $z_u^T \cdot z_v$

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log \left(\frac{\exp(\mathbf{z}_u^T \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^T \mathbf{z}_n)} \right)$$


- Problem in second Σ over all nodes
- Hard to find optimal solution

Negative Sampling

- Use *Negative Sampling* to approximate denominator

$$\log \left(\frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)} \right)$$

random distribution
over all nodes

$$\approx \log(\sigma(\mathbf{z}_u^\top \mathbf{z}_v)) - \sum_{i=1}^k \log(\sigma(\mathbf{z}_u^\top \mathbf{z}_{n_i})), n_i \sim P_V$$


from Leskovec et al., 2018

- Sample in proportion to node degree
- Experiment with k to impact negative prior and robustness
- No need to sample non-connected edges — same as random

Feature representation

- How to construct pair of nodes representation having node embeddings?
- Will it be more efficient than $\sigma(z_i^t \cdot z_j)$

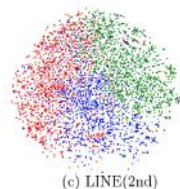
Symmetry operator	Definition
Average	$\frac{f_i(u) + f_i(v)}{2}$
Hadamard	$f_i(u) \cdot f_i(v)$
Weighted-L ₁	$ f_i(u) - f_i(v) $
Weighted-L ₂	$(f_i(u) - f_i(v))^2$
Neighbor Weighted-L ₁	$\left \frac{\sum_{w \in N(u) \cup \{u\}} f_i(w)}{ N(u) + 1} - \frac{\sum_{t \in N(v) \cup \{v\}} f_i(t)}{ N(v) + 1} \right $
Neighbor Weighted-L ₂	$\left(\frac{\sum_{w \in N(u) \cup \{u\}} f_i(w)}{ N(u) + 1} - \frac{\sum_{t \in N(v) \cup \{v\}} f_i(t)}{ N(v) + 1} \right)^2$

DOI: [10.7717/peerj-cs.172/table-2](https://doi.org/10.7717/peerj-cs.172/table-2)

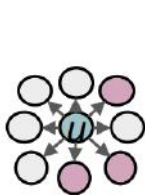
- Finite unbiased random walks (DeepWalk)
- 1-hops & 2-hops for half of the embedding, arbitrary random walks for the second half of the embedding (LINE)
- Diffusion for sampling (Diff2Vec)
- Biased random walks combining BFS and DFS (Node2Vec)

DeepWalk & LINE

- DeepWalk: Unbiased random walks with fixed length and number
- LINE: Combination of first-order and second-order proximity aggregation via concatenation



- BFS samples local neighborhood, DFS goes for global features



BFS:

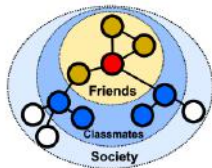


DFS:

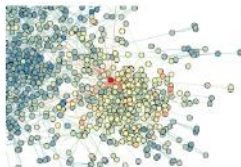
- Two parameters p and q to control sampling
- Second order Markov process

Grarep & Walklets

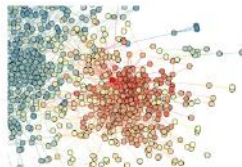
- Grarep: Approximate normalized A^k efficiently
- Walklets: approximate attention over k -distance neighbors



(a) A student (in red) is a member of several increasing larger social communities.

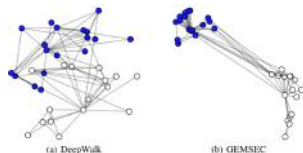


(b) WALKLETS Fine Representation



(c) WALKLETS Coarse Representation

- GEMSEC model cluster information adding regularization term over community labels given number of classes
- parameter γ balance cluster error and cluster structure given model hyper-parameters



$$\mathcal{L} = \underbrace{\sum_{v \in V} \left[\ln \left(\sum_{u \in V} \exp(f(v) \cdot f(u)) \right) - \sum_{n_i \in N_S(v)} f(n_i) \cdot f(v) \right]}_{\text{Embedding cost}} + \underbrace{\gamma \cdot \sum_{v \in V} \min_{c \in C} \|f(v) - \mu_c\|_2}_{\text{Clustering cost}} \quad (5)$$

VERSE for Structural Similarity

- VERSE uses conditional probability not on the embedding, but on the similarity rank.
- $O(|E|)$ performance

$$\text{sim}_E(v, \cdot) = \frac{\exp(W_v W^\top)}{\sum_{i=1}^n \exp(W_v \cdot W_i)}$$

$$\mathcal{L} = - \sum_{v \in V} \text{sim}_G(v, \cdot) \log(\text{sim}_E(v, \cdot))$$

$$\mathcal{L}_{NCE} = \sum_{\substack{u \sim \mathcal{P} \\ v \sim \text{sim}_G(u, \cdot)}} \left[\log \Pr_W(D = 1 | \text{sim}_E(u, v)) + \right.$$

$$\left. s \mathbb{E}_{\tilde{v} \sim Q(u)} \log \Pr_W(D = 0 | \text{sim}_E(u, \tilde{v})) \right]$$

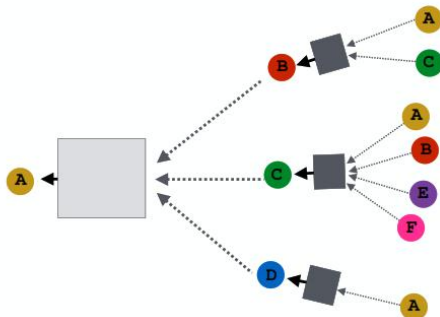
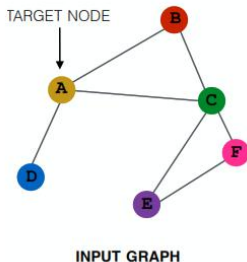
GNN

Graph Neural Network: Setting

- We have a graph $G(V, E)$ defined by adjacency matrix A and feature matrix $X \in \mathbb{R}^{f, |V|}$
- Confirmed relation between closeness of feature space and graph structure
- Non-graph features are vectorized separately (images, texts, one-hot encoding for labels, numeric features)

Graph Neural Network: Idea

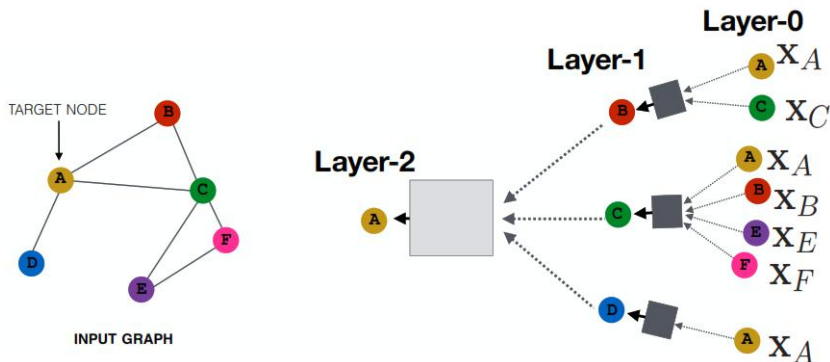
- Assign weights only to information obtained from neighbors
- Include node itself via loop with trainable weight
- Each node generate its own computational graph



from Leskovec et al., 2018

Graph Neural Network: Layer structure

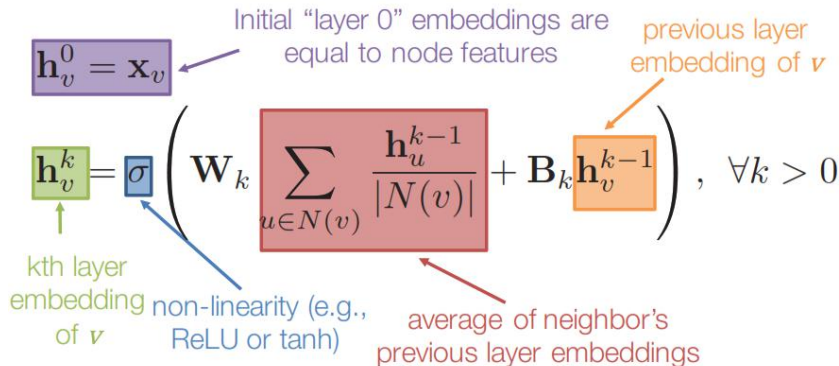
- Each aggregation defines new layer
- Zero-level embedding is non-graph feature
- Arbitrary depth but remember on “law of six handshakes”



from Leskovec et al., 2018

Graph Neural Network: Basic Approach

- Aggregation over weighted sum of neighbor input and node itself under non-linearity
- Use simple neural network construction



from Leskovec et al., 2018

GCN

Graph Convolutional Network

- Aggregation over shared weights between node and its neighbors
- Normalization to stabilize training for high-degree nodes

Basic Neighborhood Aggregation

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right)$$

VS.

GCN Neighborhood Aggregation

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)| |N(v)|}} \right)$$

same matrix for self and
neighbor embeddings

per-neighbor normalization

Graph Convolutional Network

- Efficient batch computation in matrix form
- Obtained $O(|E|)$ complexity (see pyG, DGL libraries)

$$\mathbf{H}^{(k+1)} = \sigma \left(\mathbf{D}^{-\frac{1}{2}} \tilde{\mathbf{A}} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(k)} \mathbf{W}_k \right)$$

$$\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$$

$$\mathbf{D}_{ii} = \sum_j \mathbf{A}_{i,j}$$

from Leskovec et al., 2018

GAT

Graph ATtention Network

- Not all the neighbors are equal

$$e_{ij} = a(\mathbf{W}\vec{h}_i, \mathbf{W}\vec{h}_j)$$

$$\alpha_{ij} = \text{softmax}_j(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}_i} \exp(e_{ik})}$$

$$\alpha_{ij} = \frac{\exp\left(\text{LeakyReLU}\left(\vec{a}^T[\mathbf{W}\vec{h}_i \parallel \mathbf{W}\vec{h}_j]\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\vec{a}^T[\mathbf{W}\vec{h}_i \parallel \mathbf{W}\vec{h}_k]\right)\right)}$$

\parallel is the concatenation operation.

$$\vec{h}'_i = \sigma\left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W}\vec{h}_j\right)$$

Graph ATtention Network

- Multi-head attention works better like in different convolution filters
- Final layer require pooling instead of concatenation

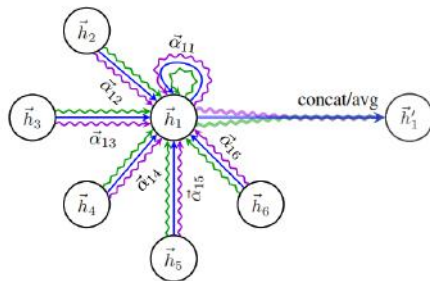
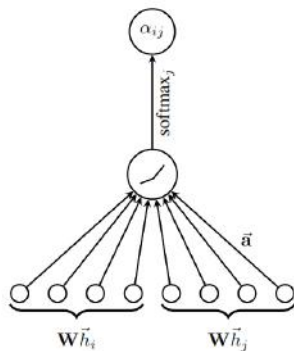
$$\vec{h}'_i = \sigma \left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W} \vec{h}_j \right)$$

$$\vec{h}'_i = \parallel_{k=1}^K \sigma \left(\sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j \right)$$

$$\vec{h}'_i = \sigma \left(\frac{1}{K} \sum_{k=1}^K \sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j \right)$$

Graph ATtention Network

- Feature aggregation via attention over learned weights
- Different patterns for the same structure



from Bengio et al., 2018

GraphSAGE

GraphSAGE: Feature Pyramid

- Vary feature space across layers
- Aggregate from neighbors and concatenate with self-representation

Simple neighborhood aggregation:

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right)$$

GraphSAGE:

concatenate self embedding and neighbor embedding

$$\mathbf{h}_v^k = \sigma \left(\left[\mathbf{W}_k \cdot \text{AGG} \left(\{ \mathbf{h}_u^{k-1}, \forall u \in N(v) \} \right), \mathbf{B}_k \mathbf{h}_v^{k-1} \right] \right)$$

generalized aggregation

Mean:

$$\text{AGG} = \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|}$$

Pool

$$\text{AGG} = \gamma(\{\mathbf{Q}\mathbf{h}_u^{k-1}, \forall u \in N(v)\})$$

element-wise mean/max

LSTM:

- Apply LSTM to random permutation of neighbors.

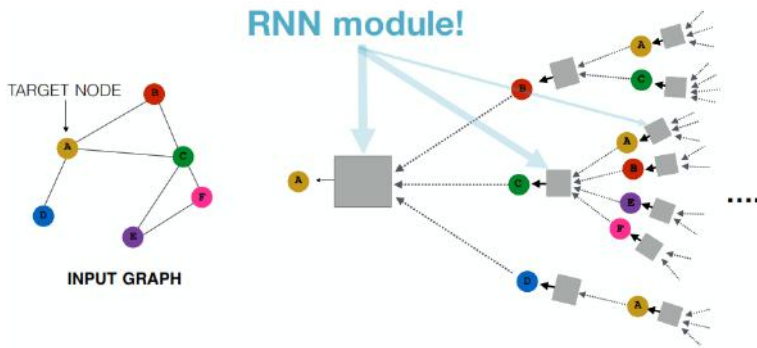
$$\text{AGG} = \text{LSTM}([\mathbf{h}_u^{k-1}, \forall u \in \pi(N(v))])$$

from Leskovec et al., 2018

How to fight dimension curse

Model Depth

- Usually 2-3 layers for GCN / GraphSAGE
- More layers make method global
- Computation graph exceed memory limits
- Overfitting, vanishing gradient



from Leskovec et al., 2018

- Use recurrent model with shared weights across all the layers, support any depth

1. Get “message” from neighbors at step k :

$$\mathbf{m}_v^k = \mathbf{W} \sum_{u \in N(v)} \mathbf{h}_u^{k-1}$$

← aggregation function does not depend on k

2. Update node “state” using Gated Recurrent Unit (GRU). New node state depends on the old state and the message from neighbors:

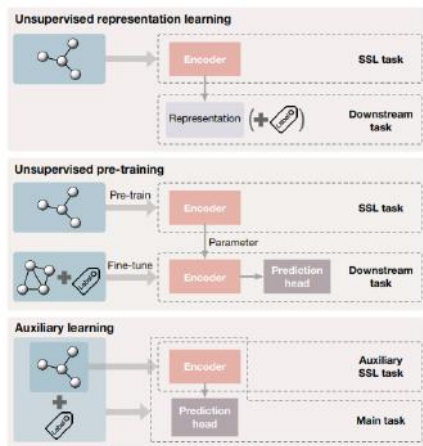
$$\mathbf{h}_v^k = \text{GRU}(\mathbf{h}_v^{k-1}, \mathbf{m}_v^k)$$

from Leskovec et al., 2018

Already too complex?
What comes next?

Self-supervised Learning

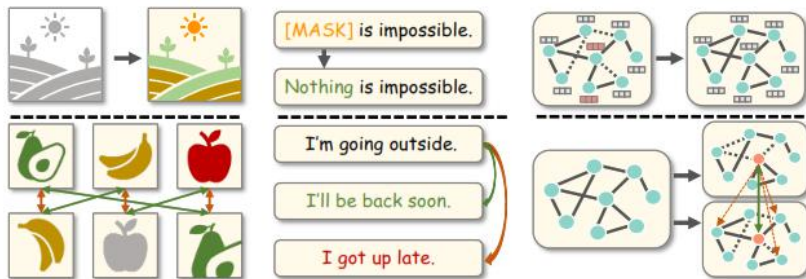
- Train data representation before using it in downstream tasks



from Ji et al., 2021

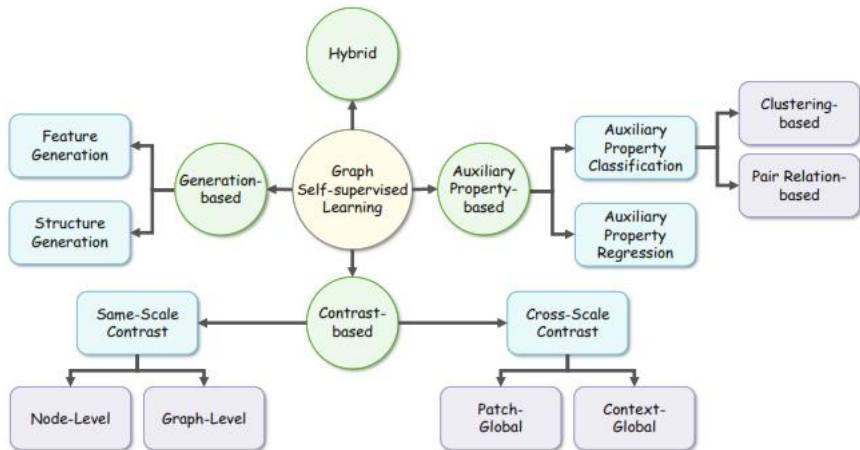
Self-supervised Learning

- Various data type lead to different data augmentation strategies



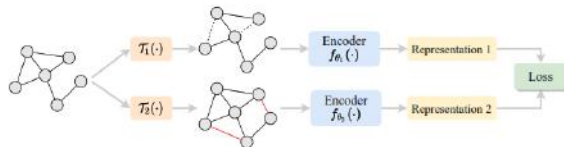
from Zhou et al., 2021

SSL for Graphs



from Zhou et al., 2021

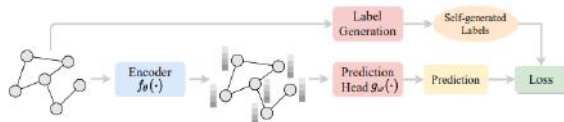
SSL for Graphs



(a) Contrastive Method



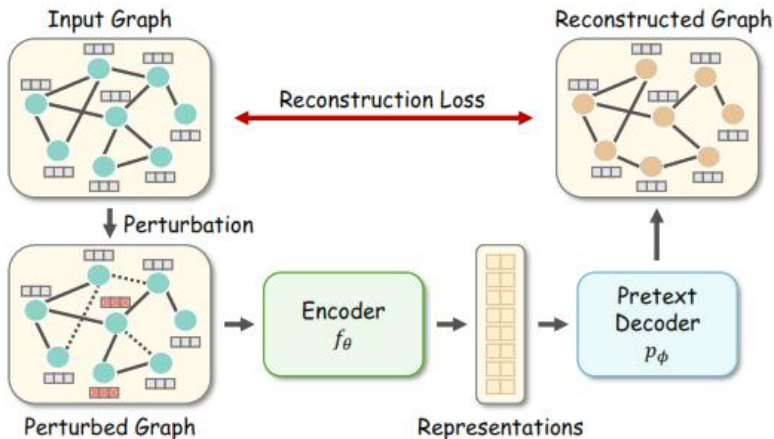
(b) Generative Method



(c) Predictive Method

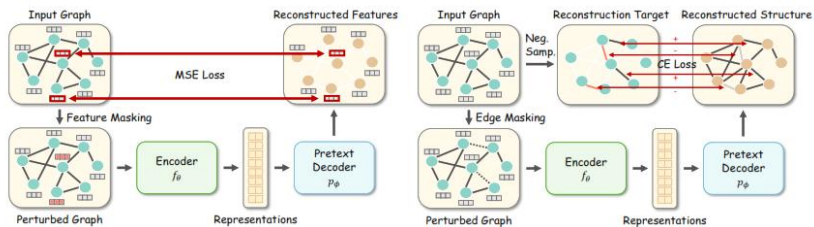
1. Generative SSL

- Use generative models trained via reconstruction loss
- Model input is generated by (optional) graph perturbation. Decoder tries to recover the original graph



1. Generative SSL methods

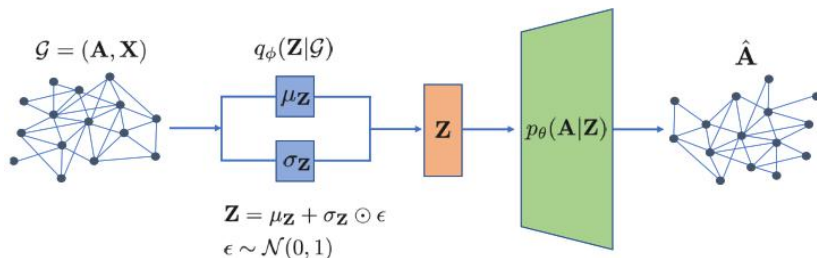
- Attribute or structure reconstruction/prediction



from Zhou et al., 2021

1. Generative SSL: Variational Graph Auto-Encoder

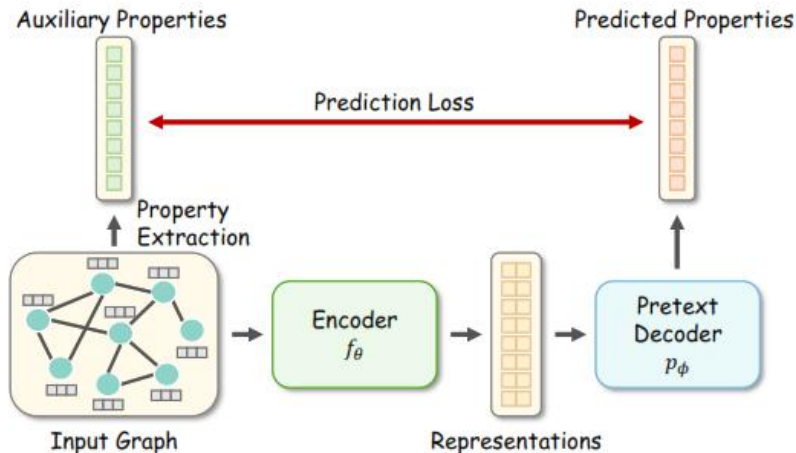
- VGAE is a variational counter-part of GAE
- Instead of direct reconstruction in GAE via LPP, VGAE learns a conditional distribution over adjacency matrices generated from embeddings



from Hamilton W.L., 2020

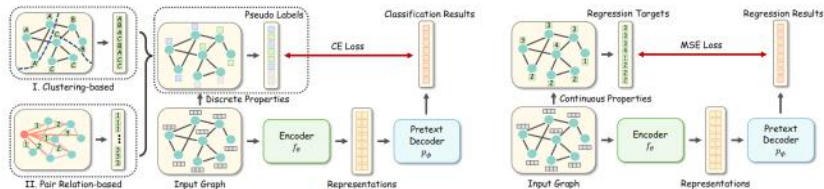
2. Predictive SSL

- Auxiliary properties are extracted from graphs freely
- Classification or regression is trained via CE/MSE loss



2. Predictive SSL methods

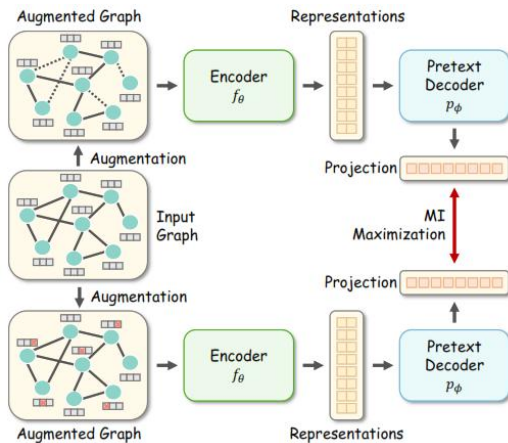
- Clustering/similarity/centrality/prestige based labels for CE/MSE loss



from Zhou et al., 2021

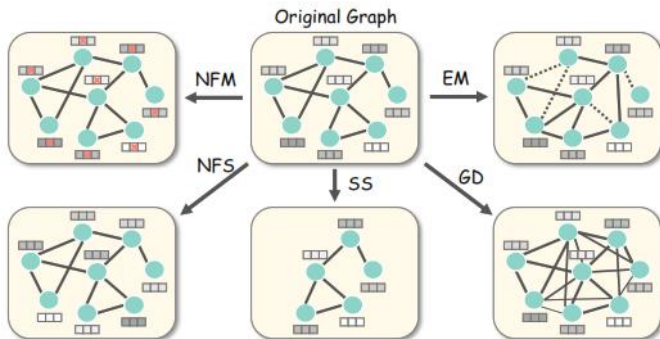
3. Contrastive SSL

- Two different augmented views are constructed from the original graph.
- Model is trained by maximizing the Mutual Information (MI) between two views.



3. Contrastive SSL: Graph augmentations

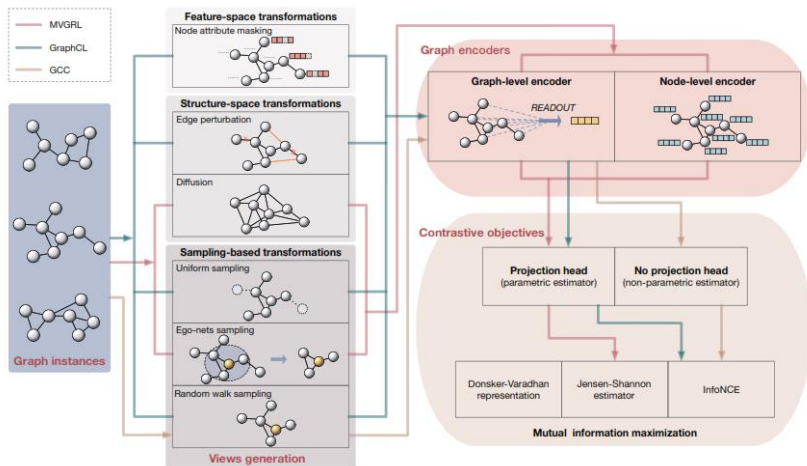
- Typical graph augmentations: Node Feature Masking (NFM), Node Feature Shuffle (NFS), Edge Modification (EM), Graph Diffusion (GD), and Subgraph Sampling (SS).



from Zhou et al., 2021

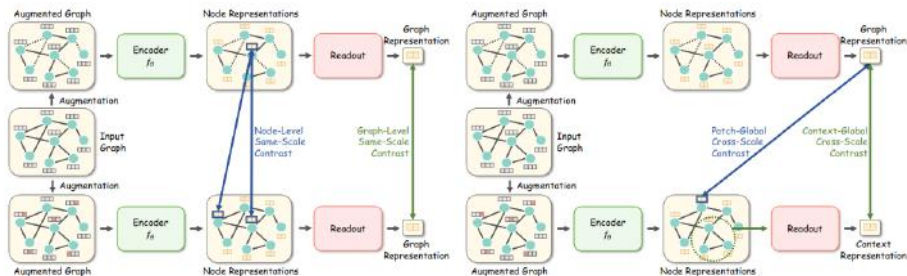
3. Contrastive SSL: Graph augmentations

- Different combinations of augmentations



3. Contrastive SSL methods

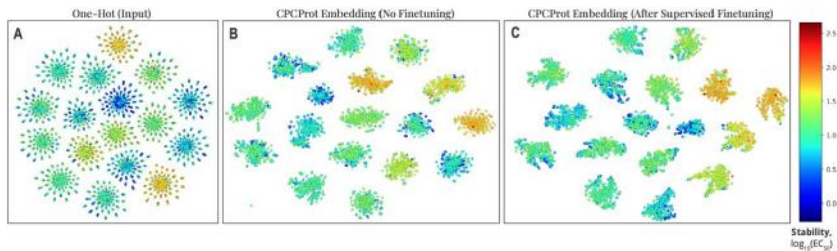
- Hierarchical augmentations on different graph levels



from Zhou et al., 2021

3. Contrastive SSL methods: SSL for Protein sequences

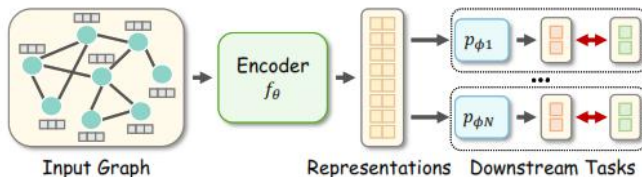
- Encoding protein sequences one can use contrastive learning to group common sequences by either, sequential patterns of predicted labels from downstream tasks



from Moses et al., 2020

4. Hybrid SSL methods

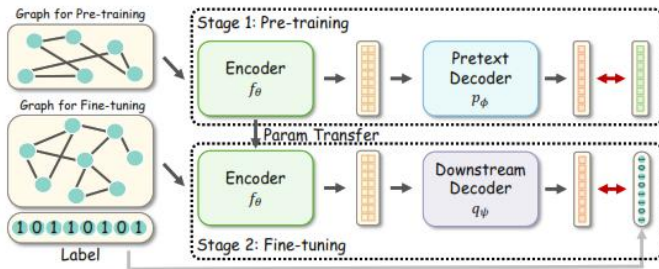
- Multiple pretext tasks are designed to train the model together in a multi-task learning manner.



from Zhou et al., 2021

Training SSL models: Pretrain and Finetune

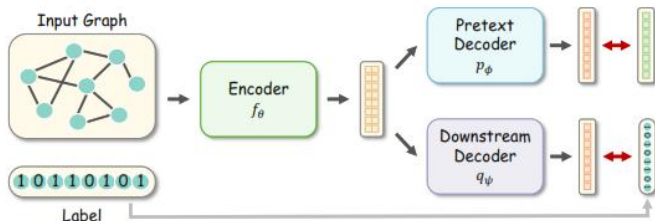
- Encoder is pretrained with pretext tasks in an unsupervised manner.
- Pretrained parameters are leveraged as the initial parameters in the fine-tuning phase for downstream tasks.



from Zhou et al., 2021

Training SSL models: Joint Learning

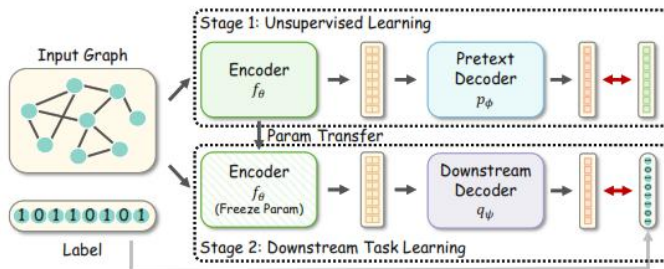
- The model is trained with pretext and downstream tasks simultaneously in a multi-task learning setting.



from Zhou et al., 2021

Training SSL models: Unsupervised Learning

- Training encoder with pretext tasks.
- Use obtained representations to learn the downstream decoder in the second phase.

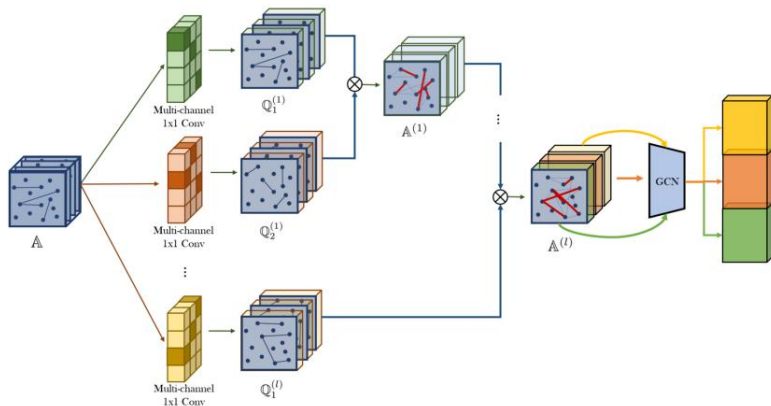


from Zhou et al., 2021

What about transformers?

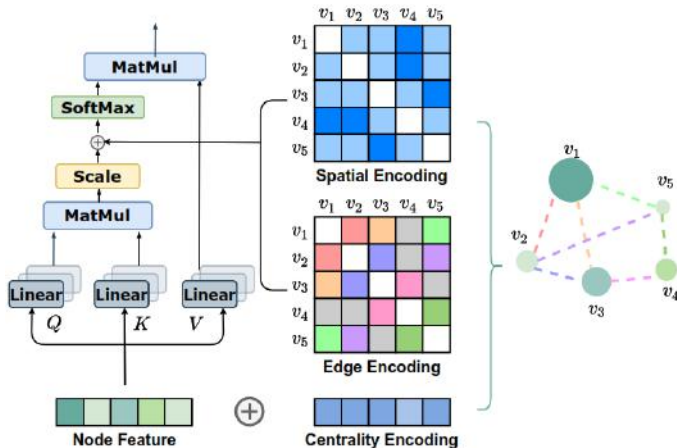
Graph Transformers

- Heterogeneous graph aggregation
- SSL Pretraining



Graph Transformers

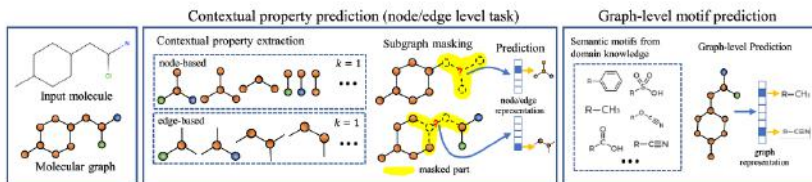
- Node and edge encodings
- Positional and similarity/distance encodings



from Liu et al., 2021

SSL for Graph Transformers

- Large molecules dataset
- Physics-/Chemistry- inspired augmentations make more sense than arbitrary augmentations
- Road towards Graph SSL Transformers!



from Huang et al., 2020

- Liu, Y., Pan, S., Jin, M., Zhou, C., Xia, F. and Yu, P.S., 2021. Graph self-supervised learning: A survey. arXiv preprint arXiv:2103.00111.
- Xie, Y., Xu, Z., Zhang, J., Wang, Z. and Ji, S., 2021. Self-supervised learning of graph neural networks: A unified review. arXiv preprint arXiv:2102.10757.
- Wu, L., Lin, H., Gao, Z., Tan, C. and Li, S., 2021. Self-supervised on Graphs: Contrastive, Generative, or Predictive. arXiv preprint arXiv:2105.07342.
- Hamilton, W.L., 2020. Graph representation learning. Synthesis Lectures on Artificial Intelligence and Machine Learning, 14(3), pp.1-159.

- Lu, A.X., Zhang, H., Ghassemi, M. and Moses, A.M., 2020. Self-supervised contrastive learning of protein representations by mutual information maximization. *BioRxiv*.
- Yun, S., Jeong, M., Kim, R., Kang, J. and Kim, H.J., 2019. Graph transformer networks. *Advances in Neural Information Processing Systems*, 32, pp.11983-11993.
- Ying, C., Cai, T., Luo, S., Zheng, S., Ke, G., He, D., Shen, Y. and Liu, T.Y., 2021. Do Transformers Really Perform Bad for Graph Representation?. *arXiv preprint arXiv:2106.05234*.
- Rong, Y., Bian, Y., Xu, T., Xie, W., Wei, Y., Huang, W. and Huang, J., 2020. Self-supervised graph transformer on large-scale molecular data. *arXiv preprint arXiv:2007.02835*.