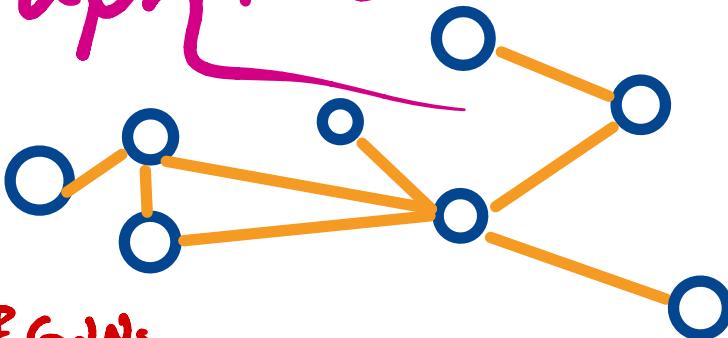


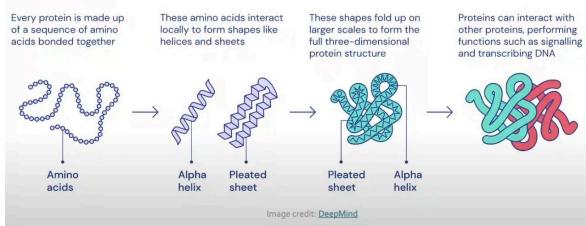
# Graph Neural Nets



## tasks of GNNs

- ① **Node level Prediction:** structure and position of a node
- ② **Edge level Prediction :** drug side effects or graph-based recom systems
- ③ **Graph level prediction:** of an entire graph or subgraphs - drug discovery, physics simulations

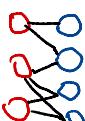
A protein chain acquires its native 3D structure



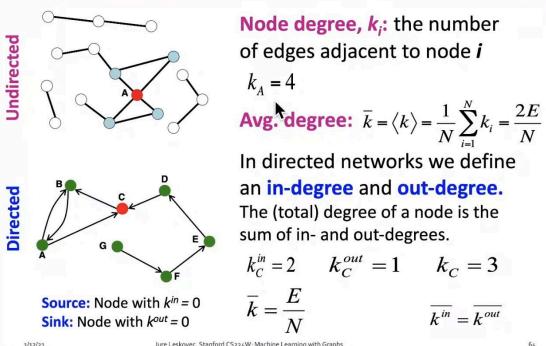
AlphaFolds key idea: spatial Graph  
nodes: amino acids  
edges: proximity between amino acids

## Bipartite Graph

a graph whose nodes can be divided into two sets  $V$  and  $U$  such that  $V$  and  $U$  are two independent sets and only interact with nodes of the other set.



## Node Degrees



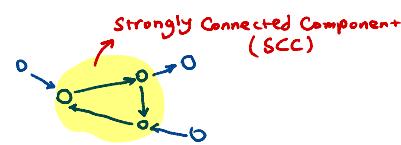
\* adjacency matrix is extremely sparse.

## Connectivity

An **undirected connected graph** is one that has at least one path between each two nodes.

in **directed** graphs we have:

- **Strongly Connected:** has a path from each node to every other and vice versa.
- **Weakly Connected:** is connected if we ignore the direction of edges.



# Traditional ML Pipeline

train an ML model  $\rightarrow$  feed new graph/node  $\rightarrow$  obtain features  $\rightarrow$  predict

## Nodes

given  $G = (V, E)$ , learn function  $f: V \rightarrow \mathbb{R}$

Characterizing the properties of a node

### Degree

$\downarrow$  makes nodes of the same degree indistinguishable.

### Centralities

while degree only counts the neighboring nodes, Centrality takes their importance into account.

### Eigenvector Centrality

the more important neighbors are the higher importance of the node.

$$Cv = \frac{1}{\lambda} \sum_{u \in \text{neighbors}} Cu \iff \lambda C = AC$$

$\lambda$ : adjacency matrix  
 $C$ : centrality vector

### Betweenness Centrality

high value if it lies on many shortest paths between other nodes. (transition hub)

### Closeness Centrality

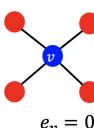
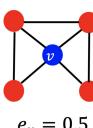
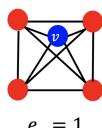
high if a node has small shortest path length to all other nodes

$$c_v := \sum_{u \neq v} \frac{1}{|\{\text{shortest paths between } u, v\}|}$$

### Clustering Coefficient

measure how connected v's neighbors are.

$$e_v := \frac{1}{\binom{k_v}{2}} |\{\text{edges among } N(v)\}| \in [0, 1]$$



### A Note on Eigenvalue Centrality (why eigenvalues?)

eigenvalue centrality is the weighted sum of neighbor centrality. This is a recursive definition. It's also similar to the definition of eigenvalues.

$$\lambda v = Av$$

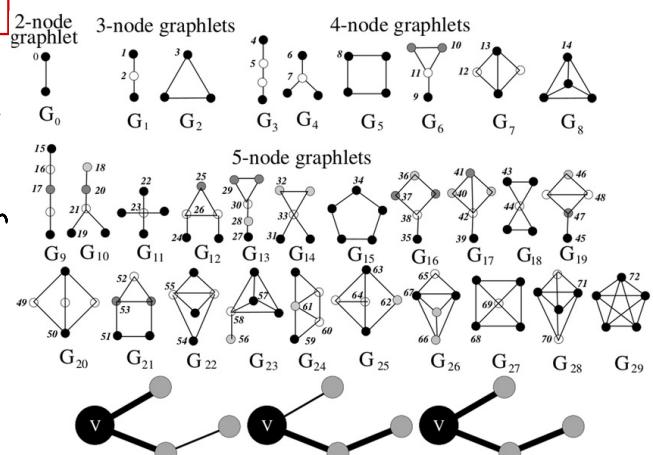
For each component of  $v$ :

$$\sum_j A_{ij} v_j = \lambda v_i$$

this means that  $v_i$  is the weighted sum of the components of its neighbor  $j$  weighted by  $A_{ij}$ . This is in definition similar to the centrality we are looking for.

### Graphlets

Small Subgraphs that describe  $v$ 's neighborhood. This can be encoded as a vector of degrees in the graphlet known as graphlet Degree Vector (GDV)



Orbit	0	1	2	3	4	5	6	...	72
GDV(v)	2	1	1	0	0	1	0	...	0

# Edge (link-level)

based on existing links, rank all node pairs and select top k pairs. (social network)

1. links missing at random

2. Links over time

the ranking depends on our scoring system.

## Local Neighborhood overlap

number of common neighbors of the two nodes.

$$\text{Jaccard's Coefficient: } \frac{|N(u) \cap N(v)|}{|N(u) \cup N(v)|}$$

if two nodes aren't connected this would always be 0.

## Global Neighborhood overlap

### Theorem

the number of paths of length  $N$  between two nodes is  $P^{(n)}$  or powers of the adjacency matrix.

We use Katz index for GNO

$$C := \sum_{i=1}^{\infty} \beta^i A^i = (I - \beta A)^{-1} - I$$

$C$  sums over paths of len 1 up to  $\infty$ ,  
 $\beta$  is the discount factor to give lower importance to longer paths.

# Graph-level

Kernel  $K(G, G') \in \mathbb{R}$  measures similarity b/w data

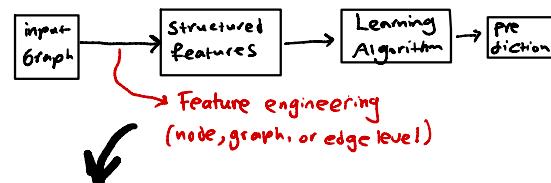
Kernel must be positive Semidefinite

a type of Kernel is Graphlet kernel and there are other types as well

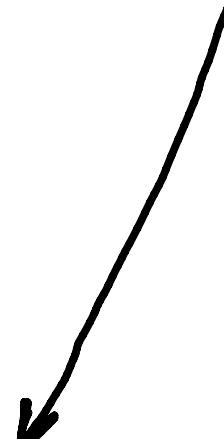
## Some Book Notes

At first glance, node classification appears to be a straightforward variation of standard supervised classification, but there are in fact important differences. The most important difference is that the nodes in a graph are not *independent and identically distributed* (i.i.d.). Usually, when we build supervised machine learning models we assume that each datapoint is statistically *independent* from all the other datapoints; otherwise, we might need to model the dependencies between all our input points. We also assume that the datapoints are *identically distributed*; otherwise, we have no way of guaranteeing that our model will generalize to new datapoints. Node classification completely breaks this i.i.d. assumption. Rather than modeling a set of i.i.d. datapoints, we are instead modeling an interconnected set of nodes.

So in summary, traditional ML for graphs looks like this:



Can we do this automatically?



# Node Embeddings

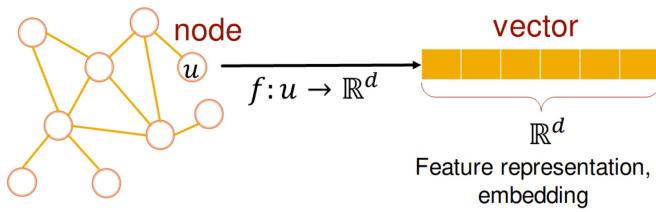


Figure 2.1: Node Level Embeddings of graphs map each node to a vector

Similarity in the embedding space must approximate similarity in the graph. this defines our goal:  
this is a pair-wise decoder

$$\text{similarity}(u, v) \simeq \text{Dec}(\text{Enc}(u), \text{Enc}(v))$$

$$\text{similarity}(u, v) \simeq z_u \cdot z_v$$

the simplest approach is to use Shallow encoder where  $Z$  is a look-up table of embeddings for each node.

How do we define note Similarity?  
this is mostly what separates these algorithms.

## Random Walk

given a graph and a starting point, we visit L number of nodes. then

$$Z_u^T Z_v \approx \begin{matrix} \text{Probability that } u \text{ and } v \\ \text{co-occur on a random walk} \end{matrix}$$

the idea: if random walk from  $u$  visits  $v$  with a high probability,  $u$  and  $v$  are similar.

goal: learn  $f: u \rightarrow \mathbb{R}^d$ :  $f(u) = z_u$

objective:

$$\max_f \sum_{u \in V} \log P(N_R(u) | z_u)$$

$N_R$  is the neighborhood of  $u$  by strategy  $R$

We want to learn representations that predict which nodes to find in  $R$ .

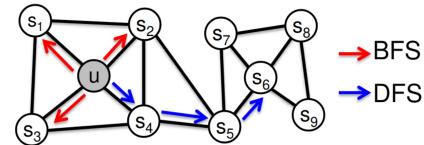
$$P(v|z_u) := \frac{\exp(z_u \cdot z_v)}{\sum_{w \in V} \exp(z_u \cdot z_w)}$$

However, this function is expensive to evaluate. The two  $\sum_{w \in V}$  loops already give  $O(|V|^2)$  time complexity. The solution to this problem is negative sampling<sup>1</sup>, which provides the estimate

$$\log \frac{\exp(z_u \cdot z_v)}{\sum_{n \in V} \exp(z_u \cdot z_n)} \simeq \log \sigma(z_u \cdot z_v) - \sum_{i=1}^k \log \sigma(z_u \cdot z_n) \quad (n_i \sim P_v)$$

## Node2vec

Similar to Randomwalk but has a flexible notion of network neighborhood.



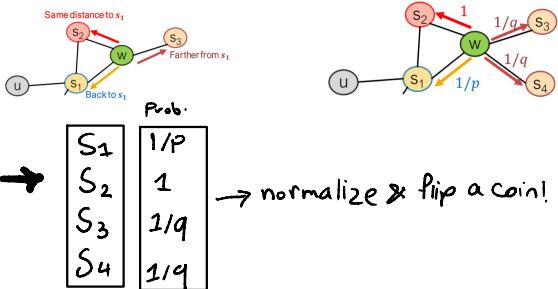
BFS: gives a Local microscopic view

DFS: gives a Global microscopic view

node2vec can extrapolate between the two

It has two hyperparams:

1. Return ( $p$ ): how likely to return to previous
2. In-out ( $q$ ): ratio of BFS vs. DFS



node2vec performs good on  
node classification. Randomwalks are  
pretty efficient as well.

Predicted probability of  $u, v$  co-occurring

Summation over all vertices

$$L(\text{Enc}) := - \sum_{u \in V} \log P(N_R(u) | z_u) = - \sum_{u \in V} \sum_{v \in N_R(u)} \log P(v | z_u)$$

Summation over nodes seen on random walk from  $u$

# Graph Embedding

a simple idea is to just average over the embeddings of nodes in a graph

another idea is to sample anonymous walks embed them, and concatenate their embeddings to get the embedding of the graph.

there are more advanced methods (later)

Page Rank is solved using power iteration

1. Initialize  $r^0$  with some value

2. Iterate  $r^{(t+1)} = M \cdot r^t$

3. Stop when  $|r^{(t+1)} - r^t| < \epsilon$

Limitations of node embeddings via random walks (or PageRank) is that we can't obtain embeddings for nodes not seen in training set.

Another limitation is these embeddings cannot capture structural similarity. So if nodes with the same neighbors are in different parts of the graph we get similar embeddings.

# Graph as a Matrix

## PageRank

the idea behind Google

Web in the old days was a collection of static webpages connected with links

"Can we rank webpages"?

webpage  $\rightarrow$  Node

links  $\rightarrow$  edge

An idea is to treat links as votes, the more incoming links, the more important that page becomes. This becomes a recursive problem.

This is solved by Stochastic Adjacency Matrix

let page  $j$  have  $d_j$  out-links,

if  $j \rightarrow i$ , then  $M_{ij} = \frac{1}{d_j}$

Columns of  $M$  sum up to 1

We define Rank Vector  $R$ , as  $r_i$  the importance score of page  $i$   $\sum_i r_i = 1$  then we have

$$R = M \cdot R$$

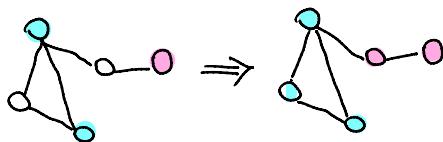


## Solution?

## Deep Representation Learning and Graph Neural Networks

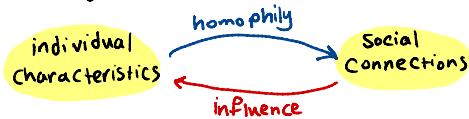
# ML with Graphs

an example is classifying nodes in a graph using labels of some of the nodes.  
(Semi-supervised)



## message Passing

Similar to PageRank, labels of a node are influenced by their neighbors.



One way is collective classification where the labels are influenced by 1st degree neighbors.  
(Markov Assumption)

Probabilistic Relational Classifier measure the weighted average of labels in the neighborhood.

Each node in a path  
listens, updates and passes a message to its neighbors.



Each node collects info from its neighbors,  
and considers the prior belief of what label it must have.

$$\begin{aligned} h_u^{(k+1)} &= \text{UPDATE}^{(k)}\left(h_u^{(k)}, \text{AGGREGATE}^{(k)}(\{h_v^{(k)}, \forall v \in \mathcal{N}(u)\})\right) \\ &= \text{UPDATE}^{(k)}\left(h_u^{(k)}, m_{\mathcal{N}(u)}^{(k)}\right), \end{aligned}$$

$$m_{\mathcal{N}(u)} = \sum_{v \in \mathcal{N}(u)} h_v,$$

$$\text{UPDATE}(h_u, m_{\mathcal{N}(u)}) = \sigma(W_{\text{self}} h_u + W_{\text{neigh}} m_{\mathcal{N}(u)}),$$

So the basic message passing can be thought of as a MLP

$$h_u^{(k)} = \sigma \left( W_{\text{self}}^{(k)} h_u^{(k-1)} + W_{\text{neigh}}^{(k)} \sum_{v \in \mathcal{N}(u)} h_v^{(k-1)} + b^{(k)} \right),$$

## Deep Learning in Graphs

a naive approach is to feed the adjacency matrix  $A$  to a feed forward network.

- $O(|V|)$  parameters
- Graph size is inherently baked into the size of the neural network
- Sensitive to node ordering



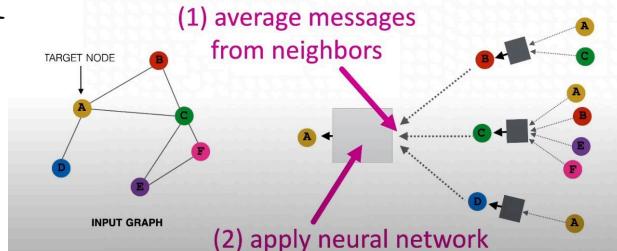
the solution is to use a permutation invariant aggregation method.

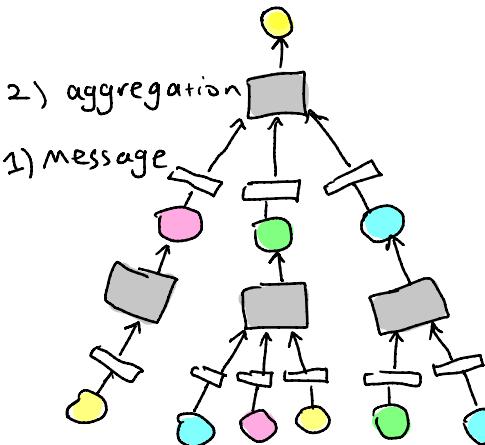
Stanford CS224W: Machine Learning with Graphs | 2021 | Lecture 6.3 - Deep Learning for Graphs

### Neighborhood Aggregation

**Basic approach:** Average information from neighbors and apply a neural network

(1) average messages from neighbors





different algorithms differ in how they aggregate messages and what messages they pass between nodes.

the idea of a single GNN layer is compress a set of message vectors into a single vector.

GNN Layer = Message + Aggregation

### 1) Message Computation

$$m_u^{(l)} = \text{MSG}^{(l)}(h_u^{(l-1)})$$

An example could be a linear transformation on the embedding of the node.

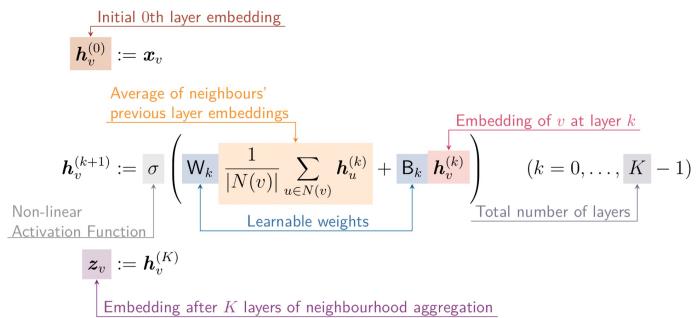
$$m_u^{(l)} = W^{(l)} h_u^{(l-1)}$$

2) Aggregation must be order-invariant as order of input nodes should not matter.  
Examples: Sum(.), Mean(.), Max(.)

but this way, when computing  $h_v^{(l)}$  the information we already have of node  $v$  gets lost! Solution?

Compute message from node  $v$  as well, using a different computation and for aggregation, Concat the two together.

$$h_v^{(l)} = \text{CONCAT}(\text{AGG}(\{m_u^{(l)}, u \in N(v)\}), m_v^{(l)})$$



## Graph Convolutional Network (GCN)

actions are

where

$$\begin{aligned} h_v^{(l)} &:= \sigma \left( \sum_{u \in N(v)} W^{(l)} \frac{1}{h_u^{(l-1)} \deg v} \right) \\ &\quad \text{Aggregation} \\ &\quad \text{Normalized by node degree} \\ m_u^{(l)} &= \frac{1}{N(v)} W^{(l)} h_u^{(l-1)} \\ h_v^{(l)} &= \sigma \left( \sum \{m_u^{(l)} : u \in N(v)\} \right) \end{aligned}$$

GCN Graphs are assumed to have self-edges that are included in the summation.

message transformation is linear and aggregation is a sum.

## Graph Attention

we keep weights of relationship between node  $v$  and its neighbors  $u$ .

$$h_v^{(l)} := \sigma \left( \sum_{u \in N(v)} \alpha_{v,u} W^{(l)} h_u^{(l-1)} \right)$$



The attention mechanism computes  $\alpha_{v,u}$ . Define the attention coefficients

$$e_{v,u} := a(W^{(l)} h_u^{(l-1)}, W^{(l)} h_v^{(l-1)})$$

Then we normalise  $e_{v,u}$  into the attention weight using softmax:

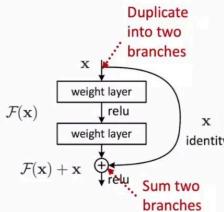
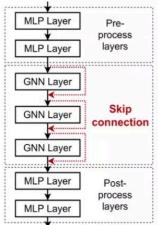
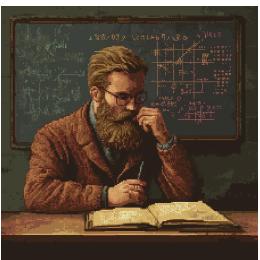
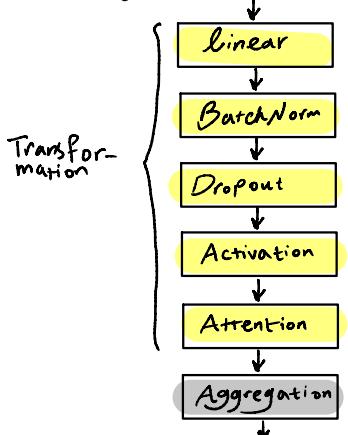
$$\alpha_{v,u} := \frac{\exp e_{v,u}}{\sum_{k \in N(v)} \exp e_{v,k}}$$

In Multi-head Attention, multiple attention scores are used and the result of each attention "head" is aggregated:

$$h_v^{(l)}[j] := \sigma \left( \sum_{u \in N(v)} \alpha_{v,u}[j] W^{(l)} h_u^{(l-1)} \right)$$

$$h_v^{(l)} := \text{Agg}(h_v^{(l)}[j] : j)$$

A GNN layer could look like this:



Idea of skip connections:  
Before adding shortcuts:  $F(x)$   
After adding shortcuts:  $F(x) + x$

# Graph Augmentation

A standard way of building GNN from a layer is to simply stack them on each other.

If many layers are stacked the problem of over-smoothing occurs where the embeddings of the nodes converge.

## Receptive Field

The set of nodes that determine the embedding of a node. As the number of GNN layers increases, the receptive field of a node covers more of the graph.

So in contrast with other domains, adding more layers to a GNN does not yield better results. We should first analyse how big we need to have the receptive field.

An option would be to use layers in a GNN that don't pass messages (like MLP).

The other option is to use skip connections.

Jure Leskovec, Stanford CS224W: Machine Learning with Graphs, <http://cs224w.stanford.edu>

50

In many cases (sparse / dense) graphs, the simple input graph may not be suitable for computation.

## Feature Augmentation

Sometimes the input graph is just adj matrix and the nodes lack any features.

Solution: add unique one-hot ids to nodes or assign constant feature to all nodes.

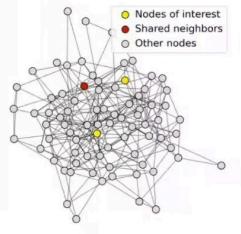
## Structural Aug

\* adding virtual (like 2-hop neighbor) nodes to a sparse graph.

\* if Graph is dense, sample node neighborhood.

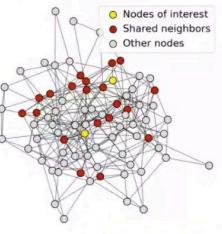
### 1-hop neighbor overlap

Only 1 node



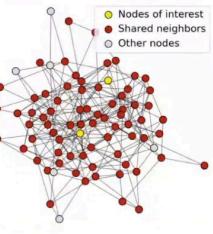
### 2-hop neighbor overlap

About 20 nodes



### 3-hop neighbor overlap

Almost all the nodes!



# Training GNN

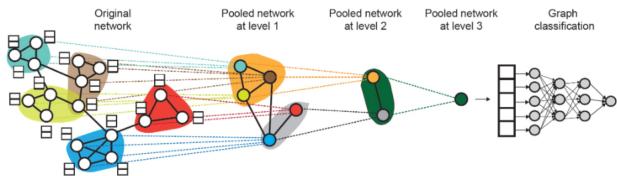


Figure 5.1: DiffPool: A Hierarchical Pooling

This is a good way of embedding entire graphs.

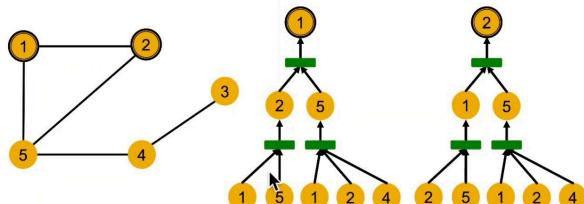
other things, such as loss functions, evaluation metrics, etc are similar to machine learning in general.

**Train/Test Split in Graphs:** Something special about Graphs is that there is information leakage meaning that since nodes are connected the splitted data in train and test are not totally independent.

## How Expressive are GNNs?

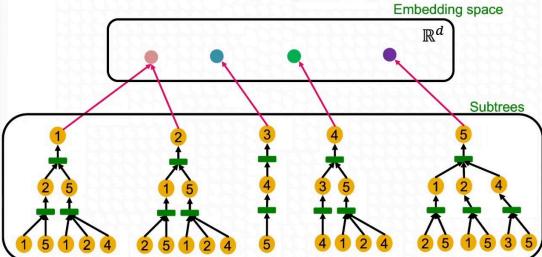
if a graph neural net has no node emb then some of the nodes would be structurally indistinguishable.

- Ex: Nodes 1 and 2's computational graphs.



without features, the computation graph of node 1 and 2 are the same.

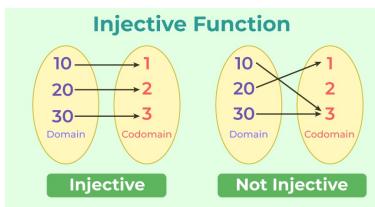
- Most expressive GNN should map subtrees to the node embeddings **injectively**.



A key factor in the expressiveness of GNNs is the aggregation function being able to distinguish different variation of neighbors.

Neighbor aggregation is a function over multi-set.

**the most expressive GNN** is the one which uses an injective function:



Theorem 6.1 (Xu et al. ICLR 2019). Any injective multi-set function can be expressed as

$$\phi \left( \sum_{x \in S} f(x) \right)$$

where  $\phi, f$  are non-linear functions.

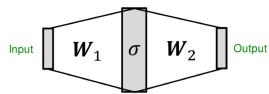


Figure 6.2: MLP with one hidden layer

Theorem 6.2 (Universal Approximation Theorem, Hornik et al., 1989). 1-hidden-layer MLP with sufficiently large hidden dimensionality and non-polynomial activation function  $\sigma$  can approximate any function to arbitrary accuracy.

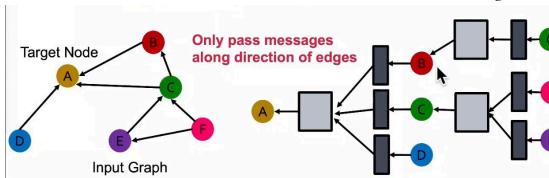
Therefore we can use the following structure to model any injective multiset function. Usually a hidden dimension of 100 to 500 is sufficient. This brings us the most expressive GNN: Graph Isomorphism Network (GIN).

$$\text{MLP}_\phi \left( \sum_{x \in S} \text{MLP}_f(x) \right)$$

GIN uses the two-MLP structure above and it is the most expressive message passing GNN.

# Relational GCN

A graph where every node and relation is labeled by a type (bio medical or event graphs)



The aggregation is similar to max pooling

$$h_v^{(l)} := \sigma \left( \sum_{u \in N^r(v)} W_r^{(l)} h_u^{(l-1)} \frac{h_u^{(l-1)}}{\deg(v)} \right)$$

Activation      Message      Aggregation

now in RGCN we have different weights for each relation type.

$$h_v^{(l)} := \sigma \left( \sum_{r \in R} \sum_{u \in N^r(v)} \frac{1}{c_{v,r}} W_r^{(l)} h_u^{(l-1)} + W_0^{(l)} h_v^{(l-1)} \right)$$

$r$ -Neighbourhood

Each relation has  $L$  matrices for each layer of the GNN:  $W_r^{(1)}, W_r^{(2)}, \dots, W_r^{(L)}$

this means rapid growth of parameters (overfitting)

2 Solutions for this:

1. Block Diagonal Mat: make the weights matrix sparse.

$$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

2. Weight-Sharing: represent  $W_r$  as a linear combination of basic transformations which are shared across all relations:  $W_r = \sum a_{rb} \cdot V_b$

basic matrix  
importance weights

## Knowledge Graphs

KGs are represented as  $(h, r, t)$  where head( $h$ ) has relation( $r$ ) with tail( $t$ ).

the goal is to make  $\text{emb}(h, r)$  close to  $\text{emb}(t)$

## TransE

If fact is true,  $h+r \approx t$

Scoring Function  $\|h+r - t\|$

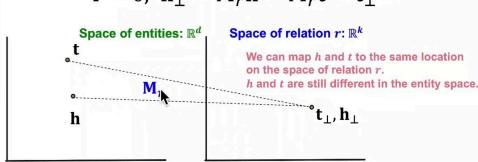
it is not able to capture symmetrical relations.



## TransR

uses  $M_r$  projection mat to transform from entity space to relation space and then perform TransE.

- TransR can model symmetric relations
- $r = 0, h_\perp = M_r h = M_r t = t_\perp$
- does Not support composition relations.  
(my mother's husband is my father)



## DistMult

Similar to TransR but uses different method to score embeddings. it can be thought of as cosine similarity between  $h \cdot r$  and  $t$ .

$$t_2 \rightarrow t_2$$

$t_1$        $h \cdot r$

$$P_r(h \cdot r, t_2) > P_r(h \cdot r, t_1)$$

Can Not model inverse relations as

$$f_{r_2}(h, r_2, t) = f_{r_1}(h, r_1, t)$$

also Can NOT model composition relations.

## ComplEx

based on DistMult but embeds into complex vector space.

$$\begin{array}{ccc} \bar{u} = a - bi & & u = a + bi \\ \bar{u} \text{ is conjugate} & & \end{array}$$

Score function:  $\text{Re}(\sum_i h_i \cdot r_i \cdot \bar{t}_i)$  takes the real part.

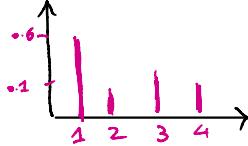
also Can NOT model composition relations.

(\* I didn't make notes for KG reasoning (next lesson))

# Generative Models for Graphs

Some Properties of real-world graphs:

1) Degree Distribution  $P(k)$  probability that a randomly chosen node has a degree of  $k$ .



2) Clustering Coefficient: how connected are the neighbors of a node. number of edges

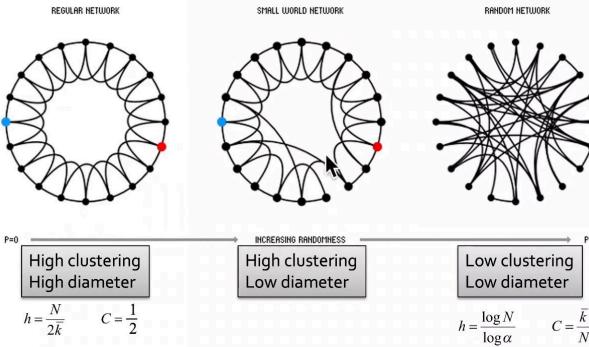
$$\binom{k}{2}$$

3) Path Length (diameter): the maximum or average of shortest path between any two nodes.

## Small World Model

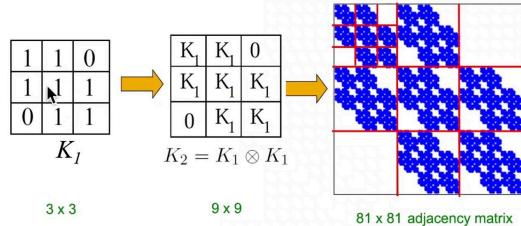
tries to have **high clustering** with **low avg. shortest path** (the two oppose each other)

- 1) Create a small lattice
- 2) rewire and add randomness



## Kronecker Model

A recursive model of network structure



the **Kronecker Product** takes every cell of  $M_1$  and multiplies it by the whole  $M_2$  matrix.

the model iteratively uses  $K_P$  on an initiator graph

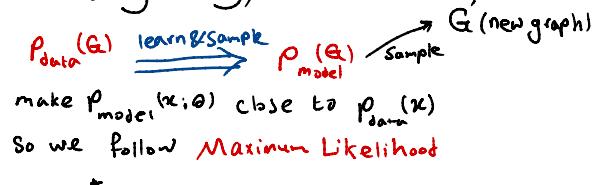
$$K_2 \otimes K_2 \otimes K_2 \otimes K_2$$

**Kronecker** and **real** graphs are similar: all you need to do is choose the right initiator matrix.

## Deep Generative Models for Graphs

task1: generate graphs similar to a given set of graphs

task2: " " that optimize given objectives (such as drug discovery)



So we follow **Maximum Likelihood**

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \sum_{x \sim p_{\text{data}}} \log P_{\text{model}}(x | \theta)$$

Idea: **Chain Rule**. we want to model a complex distribution over graphs, we don't know how, so we break it down to smaller distributions.



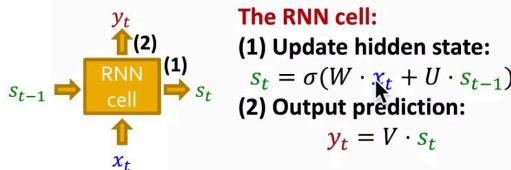
# Graph RNN

Generating a graph sequentially by adding nodes and edges.



this has become a sequencing problem → **RNN**

- $s_t$ : State of RNN after step  $t$
- $x_t$ : Input to RNN at step  $t$
- $y_t$ : Output of RNN at step  $t$
- RNN cell:  $W, U, V$ : Trainable parameters



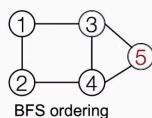
as the generation goes on, a new node can connect to any/all previous nodes. the model must remember all the nodes and edges → scale problem!

So we use BFS ordering!

## Breadth-First Search node ordering

"Recipe" to generate the left graph:

- Add node 1
- Add node 2
- Connect 2 with 1
- Add node 3
- Connect 3 with 1
- Add node 4
- Connect 4 with 3 and 2



## BFS node ordering:

- Since Node 4 doesn't connect to Node 1
- We know all Node 1's neighbors have already been traversed
- Therefore, Node 5 and the following nodes will never connect to node 1
- We only need memory of 2 "steps" rather than  $n - 1$  steps

We can also use GNNs to generate graphs (GCPN)

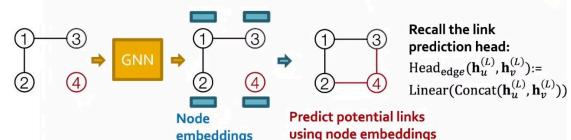
## Sequential graph generation

- GraphRNN: predict action based on **RNN hidden states**

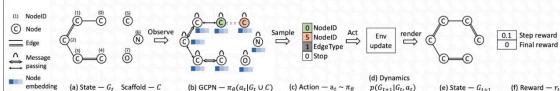


RNN hidden state captures the generated graph so far

- GCPN: predict action based on **GNN node embeddings**



Recall the link prediction head:  
 $\text{Head}_{\text{edge}}(\mathbf{h}_u^{(L)}, \mathbf{h}_v^{(L)}) := \text{Linear}(\text{Concat}(\mathbf{h}_u^{(L)}, \mathbf{h}_v^{(L)}))$



- (a) Insert nodes
- (b,c) Use GNN to predict which nodes to connect
- (d) Take an action (check chemical validity)
- (e, f) Compute reward

## Two parts:

- (1) **Supervised training:** Train policy by **imitating the action** given by real observed graphs. Use **gradient**.

■ We have covered this idea in GraphRNN

- (2) **RL training:** Train policy to **optimize rewards**. Use standard **policy gradient** algorithm.

■ Refer to any RL course, e.g., CS234



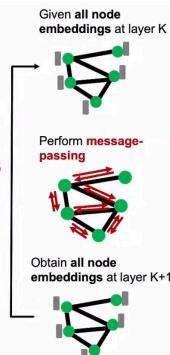
# Scaling up GNNs

If we use a mini-batch SGD training for a GNN and we select nodes that are independent of each other, we cannot effectively train GNN.

## Naïve full-batch

**implementation:** Generate embeddings of all the nodes **at the same time**.

- Load the entire graph and features
- At each GNN layer:
  - Compute embeddings of all nodes using all the node embeddings from the previous layer.
- Compute the loss
- Perform gradient descent

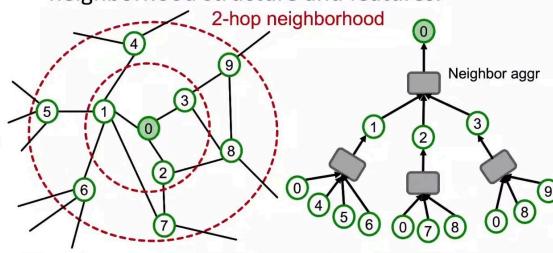


it takes so long and much computation

## GraphSAGE neighborhood Sampling

RECALL: GNN generate node embeddings via neighborhood aggregation.

- Observation:** A 2-layer GNN generates embedding of node "0" using 2-hop neighborhood structure and features.

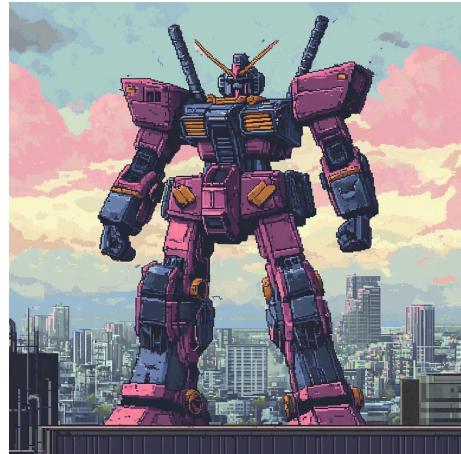


If # of layers or the neighborhood is small we can ignore much of the network now we can compute gradients in a reliable way and use SGD.

But the computation graph grows exponentially especially if we hit a hub node.



neighbor Sampling : at each hop, pick at most  $H_k$  neighbors.



X : @Hesamation